

Probabilistic Machine Learning Estimation of Ocean Mixed Layer Depth from Dense Satellite and Sparse In-Situ Observations

Dallas Foster¹, David John Gagne², and Daniel Bridger Whitt³

¹Oregon State University

²NCAR

³National Center for Atmospheric Research

November 23, 2022

Abstract

The ocean mixed layer plays an important role in subseasonal climate dynamics because it can exchange large amounts of heat with the atmosphere, and it evolves significantly on subseasonal timescales. Estimation of the subseasonal variability of the ocean mixed layer is therefore important for subseasonal to seasonal prediction and analysis. The increasing coverage of in-situ Argo ocean profile data allows for greater analysis of the aseasonal ocean mixed layer depth (MLD) variability on subseasonal and interannual timescales; however, current sampling rates are not yet sufficient to fully resolve subseasonal MLD variability. Other products, including gridded MLD estimates, require optimal interpolation, a process that often ignores information from other oceanic variables. We demonstrate how satellite observations of sea surface temperature, salinity, and height facilitate MLD estimation in a pilot study of two regions: the mid-latitude southern Indian and the eastern equatorial Pacific Oceans. We construct multiple machine learning architectures to produce weekly 1/2 degree gridded MLD anomaly fields (relative to a monthly climatology) with uncertainty estimates. We test multiple traditional and probabilistic machine learning techniques to compare both accuracy and probabilistic calibration. We find that incorporating sea surface data through a machine learning model improves the performance of MLD estimation over traditional optimal interpolation in terms of both mean prediction error and uncertainty calibration. These preliminary results provide a promising first step to greater understanding of aseasonal MLD phenomena and the relationship between the MLD and sea surface variables. Extensions to this work include global and temporal analyses of MLD.

1 **Probabilistic Machine Learning Estimation of Ocean**
2 **Mixed Layer Depth from Dense Satellite and Sparse**
3 **In-Situ Observations**

4 **Dallas Foster¹, David John Gagne II², Daniel B. Whitt²**

5 ¹Department of Mathematics, Oregon State University, Corvallis, OR 97331

6 ²National Center for Atmospheric Research, Boulder, CO, 80305

7 **Key Points:**

- 8 • Machine learning models that incorporate surface and ocean profile data improve
9 ocean MLD estimates.
- 10 • Model performance is dependent on spatial location and strength of the sub-seasonal
11 variance.
- 12 • Probabilistic sampling techniques capture uncertainty better than standard or para-
13 metric approaches.

Abstract

The ocean mixed layer plays an important role in subseasonal climate dynamics because it can exchange large amounts of heat with the atmosphere, and it evolves significantly on subseasonal timescales. Estimation of the subseasonal variability of the ocean mixed layer is therefore important for subseasonal to seasonal prediction and analysis. The increasing coverage of in-situ Argo ocean profile data allows for greater analysis of the aseasonal ocean mixed layer depth (MLD) variability on subseasonal and interannual timescales; however, current sampling rates are not yet sufficient to fully resolve subseasonal MLD variability. Other products, including gridded MLD estimates, require optimal interpolation, a process that often ignores information from other oceanic variables. We demonstrate how satellite observations of sea surface temperature, salinity, and height facilitate MLD estimation in a pilot study of two regions: the mid-latitude southern Indian and the eastern equatorial Pacific Oceans. We construct multiple machine learning architectures to produce weekly 1/2 degree gridded MLD anomaly fields (relative to a monthly climatology) with uncertainty estimates. We test multiple traditional and probabilistic machine learning techniques to compare both accuracy and probabilistic calibration. We find that incorporating sea surface data through a machine learning model improves the performance of MLD estimation over traditional optimal interpolation in terms of both mean prediction error and uncertainty calibration. These preliminary results provide a promising first step to greater understanding of aseasonal MLD phenomena and the relationship between the MLD and sea surface variables. Extensions to this work include global and temporal analyses of MLD.

Plain Language Summary

The top layer of the ocean, called the surface mixed layer, features temperature and salinity that are relatively uniform throughout its depth. The depth of this layer can vary depending on the exact location, time of year and is impacted by many physical processes. Although it is typically only a few percent of the ocean depth, the mixed layer is important because it regulates heat exchange between the deep ocean and the atmosphere, and it hosts virtually all photosynthesis that sustains ocean ecosystems. Observations of the mixed layer depth are infrequent in time and space because of the size of the ocean in comparison to the number of observing instruments. Satellite data is widely available for information about the surface of the ocean, but unfortunately there is not an exact relationship between the surface information and the mixed layer depth. In this paper, we study machine learning models' abilities to learn this relationship with the available data and to produce reasonable fine-scale estimates of the mixed layer depth. In particular, we emphasize the ability of the machine learning model to estimate how uncertain it is about its estimates.

1 Introduction

Because of the ocean surface mixed layer's role as intermediary between ocean and atmosphere, many important processes, such as water mass formation and ocean circulation (Hanawa & Talley, 2001; Stommel, 1979) and air-sea interaction (Frankignoul & Hasselmann, 1977; Kraus & Turner, 1967) are sensitive to the ocean surface mixed layer depth (MLD). While there have been several recent efforts to observe and quantify the global climatological behavior of the MLD based on the in-situ array of thousands of vertically-profiling Argo floats (Holte et al., 2017; Schmidtke et al., 2013; D. B. Whitt et al., 2019), little effort has been devoted to quantifying the subseasonal and interannual (aseasonal) variability of the MLD because the Argo array is not sufficiently large to fully resolve subseasonal MLD variability. Through this study, we take a preliminary step toward improved observational estimates of aseasonal MLD variability by investigating the relationship between MLD and sea surface salinity, temperature, and height anomalies.

64 Due largely to the increasing coverage of the Argo array (Holte et al., 2017), the
 65 MLD is increasingly well-observed globally. Despite this improvement, however, the data
 66 is insufficient to recover sub-seasonal processes on a fine grid at high frequency. Mod-
 67 ern attempts to recover variables using a hybrid data collection of in-situ and satellite
 68 data typically use optimal interpolation (Roemmich & Gilson, 2009; Guinehut et al., 2012).
 69 Our aim in this paper is to demonstrate the utility of informing MLD estimation using
 70 satellite surface data through a machine learning framework.

71 The application of machine learning to the geosciences is a rapidly growing field
 72 ((Monteleoni et al., 2013; Reichstein et al., 2019; Weyn et al., 2019; Lary et al., 2016;
 73 Irrgang et al., 2020). The machine learning approach offers a flexible, data-driven route
 74 to regression and classification tasks that has been used for parameterizations (Bolton
 75 & Zanna, 2019; Gagne et al., 2020; Rasp et al., 2018; O’Gorman & Dwyer, 2018; Gen-
 76 tine et al., 2018; Jiang et al., 2018; Brenowitz & Bretherton, 2018), forecasting (Pathak
 77 et al., 2018; McGovern et al., 2017; Ukkonen & Mäkelä, 2019; Irrgang et al., 2020; Weyn
 78 et al., 2019; Hsieh & Tang, 1998), data assimilation (R. Cintra et al., 2016; Wahle et al.,
 79 2015; R. S. Cintra & Velho, 2018), and remote sensing (Lary et al., 2016; Ouali et al.,
 80 2017). The commonality to many of these approaches and the motivation for use in this
 81 study is not only the lack of a deterministic model between the sea surface variables and
 82 the mixed layer depth, but also the possibility of an empirical model being learned from
 83 the existing data. Unfortunately, many successes in machine learning research are also
 84 in over-determined regimes, in which the amount of data is large in comparison to the
 85 number of independent parameters. Extrapolation regimes, where data are sparse in one
 86 or more dimensions, are known to be problematic because the prediction depends more
 87 heavily on the underlying assumptions of the model. This is particularly problematic in
 88 oceanography, where many unknown quantities are 2 or 3 dimensional, and data avail-
 89 ability is still relatively sparse.

90 While the study of machine learning can trace its history to Rosenblatt’s percep-
 91 tron (Rosenblatt, 1958), the implementation of early machine learning methods and archi-
 92 tectures in a data-driven way was considered computationally infeasible for moder-
 93 ate to large applications until the late 1980s with the development of the back-propagation
 94 algorithm (Rumelhart et al., 1986), which enabled training of multi-layered neural net-
 95 works. Despite advances through the nineties and early twenty-first century, the deep
 96 learning revolution did not occur until 2006 (Goodfellow et al., 2016) when an explosion
 97 of reliable training data, computing power, neural network layers, and regularization tech-
 98 niques have dramatically increased neural network accuracy. As demonstrated in Guo
 99 et al. (2017), this improvement in accuracy has also hindered the capacity of neural net-
 100 works to be well-calibrated, i.e. when forecast probabilities match the system’s true prob-
 101 abilities, and hence offer accurate representations of the underlying probability distri-
 102 butions. The ability for a neural network to be well-calibrated is of critical importance.
 103 Data Assimilation research has repeatedly shown that proper estimation of the background
 104 error covariance can improve reconstruction estimates (Valler et al., 2019). In the esti-
 105 mation of sea surface temperature or sea level anomaly, mis-quantification of atmospheric
 106 uncertainties has also been shown to cause significant and non-local errors in reanaly-
 107 sis estimates (Chaudhuri et al., 2016). Parallel developments have led to the field of prob-
 108 abilistic neural networks to address this calibration problem in machine learning.

109 The ultimate goal of probabilistic neural networks is to be able to accurately and
 110 precisely define the posterior probability distribution conditioned on the data. Using a
 111 Bayesian framework allows us to easily account for sources of error and randomness in
 112 the data, weights, or model. The gold standard for this task is often sampling from the
 113 posterior distribution using a Markov Chain Monte Carlo (MCMC) scheme (Brooks, 2011;
 114 Gelman et al., 2013), but this approach is still computationally infeasible for modern neu-
 115 ral networks. There have been several approximations and techniques developed for pro-
 116 ducing estimates of the posterior probability including the development of Bayesian Neu-
 117 ral Networks, with weight uncertainty (Neal, 1996; Blundell et al., 2015), Stochastic Gra-
 118 dient Langevin Dynamics (Welling & Teh, 2011), Variational Inference (Paisley et al.,

2012; Hoffman & Blei, 2015; Kingma et al., 2015), Probabilistic Backpropagation (Rezende et al., 2014; Hernández-Lobato & Adams, 2015), Dropout (Hinton et al., 2012; Ba & Frey, 2013; Maeda, 2014; Gal & Ghahramani, 2016; Gal et al., 2017), Variational Autoencoders (Kingma & Welling, 2014), and Deep Ensembles (Lakshminarayanan et al., 2017).

Despite the numerous techniques to inject uncertainty estimates into machine learning, the performance of any approach is still underwhelming. Recent arguments have been made that ensembles of techniques outperform any one approach (Lakshminarayanan et al., 2017; Kuleshov et al., 2018; Guo et al., 2017; Nixon et al., 2019; Dormann, 2020). Due to the complex nature of the analytical posterior distributions, lack of complete data, prohibitive cost of training, and sensitivity to the nature of the application, an understanding of which methodology is appropriate is still in its infancy. Recently there has been some research comparing popular uncertainty quantification techniques in Deep Learning (Ashukha et al., 2020; Caldeira & Nord, 2020; Labach et al., 2019; Lakshminarayanan et al., 2017). Unfortunately, there is not much research about how these methods perform in the geosciences, where probabilities are often non-Gaussian, non-trivial, non-stationary, and high-dimensional. This paper serves as a step into answering this question by testing various probabilistic machine learning methods used for high-dimensional data with both Gaussian and non-Gaussian distributions on MLD estimation, which serves as an example problem in this respect.

Our goal for this manuscript is two-fold. First, we investigate to what extent the aseasonal variability in sea surface salinity, temperature, and height are related to, and hence useful for estimating, the aseasonal variability of the MLD. In particular, we study two geographic regions, (1) the eastern equatorial Pacific Ocean from 10S-10N and 150W-120W and (2) the southern Indian Ocean from 45S-35S from 60E-120E, over the 2011-2015 time period. As detailed in section 2, these regions are useful test cases because both are characterized by at least modest subseasonal MLD variability (> 10 m subseasonal standard deviations), but the magnitudes of subseasonal variability, the climatological annual cycle, and interannual variability all differ substantially (D. B. Whitt et al., 2019). Thus, the two regions reflect useful and distinct test cases for evaluating machine learning model performance. We perform this analysis by training a series of neural network architectures to produce gridded MLD estimates using surface variables as inputs and evaluate model performance using the Argo profiles. We compare the machine learning approaches, which only use surface values as inputs, to the traditional optimal-interpolation technique that estimates using the actual MLD values from the in-situ Argo profiles. The differences in performance between the machine learning methods and optimal-interpolation schemes will reveal the extent to which the sea surface variables are useful in predicting the MLD.

Second, we focus on understanding the probability distribution of the MLD that is learned by the neural network. As a first step, we evaluate how well calibrated the neural network estimates are and what spatial and temporal patterns are revealed through sampling these distributions. We choose three probabilistic machine learning methods that cover two distinct types of uncertainty quantification: parameterization- and sampling-based methods. By evaluating these methods, we aim to understand the appropriateness of a Gaussian distribution to the data and the ability for sampling machine learning methods in exploring the posterior distribution. Finally, we compare the machine learning uncertainty quantification against uncertainty estimates from the optimal-interpolation approach. As before, this last comparison will reveal the extent to which the sea surface variables inform us about the uncertainty in the MLD.

These methods are certainly not exhaustive and so this paper is a first step to a better understanding of the aseasonal MLD variability and how machine learning can be used as a tool in this investigation. The outline of the body of the paper is as follows: first, in section 2 we detail the data and describe the data processing and methodology; second, in section 3 we describe the mathematical framework and relevant machine learn-

172 ing architectures that we implement; lastly, in section 4 we explain and detail the ex-
 173 periments and results.

174 2 Data

175 2.1 Salinity

176 Sea surface salinity data is the optimally-interpolated analysis of Melnichenko et
 177 al. (2016), which is an optimal interpolation of observations from the Aquarius satellite
 178 and uses corrections to minimize bias relative to in-situ data. The data exists on a $\frac{1}{2}$ de-
 179 gree, weekly grid spanning roughly 2011-2015 (200 weeks). A random 150 week sample
 180 constitutes the training data, with the remaining being used for testing and validation.
 181 This grid is the coarsest of all the variables and thus will form the basis that we inter-
 182 polate and re-sample the other data onto. To calculate an estimate of the climatology,
 183 we calculate monthly means using only the training data, taking a 4 week boxcar mov-
 184 ing average, binning data into months and averaging over the bins.

185 2.2 Temperature

186 Sea surface temperature data comes from the GHRSSST Level 4 Global Foundation
 187 Sea Surface Temperature analysis dataset (Remote Sensing Systems, 2017). This dataset
 188 uses Optimal Interpolation (OI) from several microwave sensors. The data exists on a
 189 $\frac{1}{4}$ degree, daily grid spanning roughly 2001-2018. To calculate an estimate of the clima-
 190 tology, we set aside the years 2011-2015 and calculate a 4 week boxcar moving average
 191 on the remaining data. From the smoothed data, we take bins according to each month
 192 and average over the bins, resulting in an approximate monthly climatology. To calcu-
 193 late anomalies, we bin the 2011-2015 data into months and subtract the monthly clima-
 194 tology. Then, to be able to compare to the salinity dataset, we up-sample from the daily
 195 values to weekly data and optimally interpolate onto a $\frac{1}{2}$ degree grid.

196 2.3 Height Anomaly

197 Sea surface height anomaly data comes from the MEaSUREs Gridded Sea Surface
 198 Height Anomalies dataset (Zlotnicki et al., 2019). The data exists on a $\frac{1}{6}$ degree, 5-day
 199 grid spanning roughly 1992-2019. We do not calculate and remove climatologies from
 200 this data set. To be able to compare to the salinity dataset, we up-sample from the 5-
 201 day values to weekly data and optimally interpolate onto a $\frac{1}{2}$ degree grid.

202 2.4 Mixed Layer Depth

203 Argo data is available through Cabanes et al. (2013). The MLD is defined for about
 204 1.5 million profiles of temperature and salinity that pass quality controls in the time span
 205 from 2000-2017 (D. B. Whitt et al., 2019; D. Whitt et al., 2020).

206 To calculate an estimate of the climatology from the individual MLD measurements,
 207 we take the years 2002-2010, and 2016-2017, bin the data into 2° latitude and 4° lon-
 208 gitude bins, re-sample onto a daily grid and take four week moving averages in each bin.
 209 This smoothed data is then grouped into months. Both an average and standard devi-
 210 ation are calculated in order to compute the mean and standard deviation of the monthly
 211 climatology in each bin.¹ Anomalies are created by taking each profile from the with-
 212 held 2011-2015 Argo data and subtracting the climatology according to the profile's bin

¹ For the regions included in our studies, all bins have enough data to calculate the monthly climatol-
 ogy. There are many regions, such as some seas surrounding Indonesia, for instance, that do not have
 sufficient data.

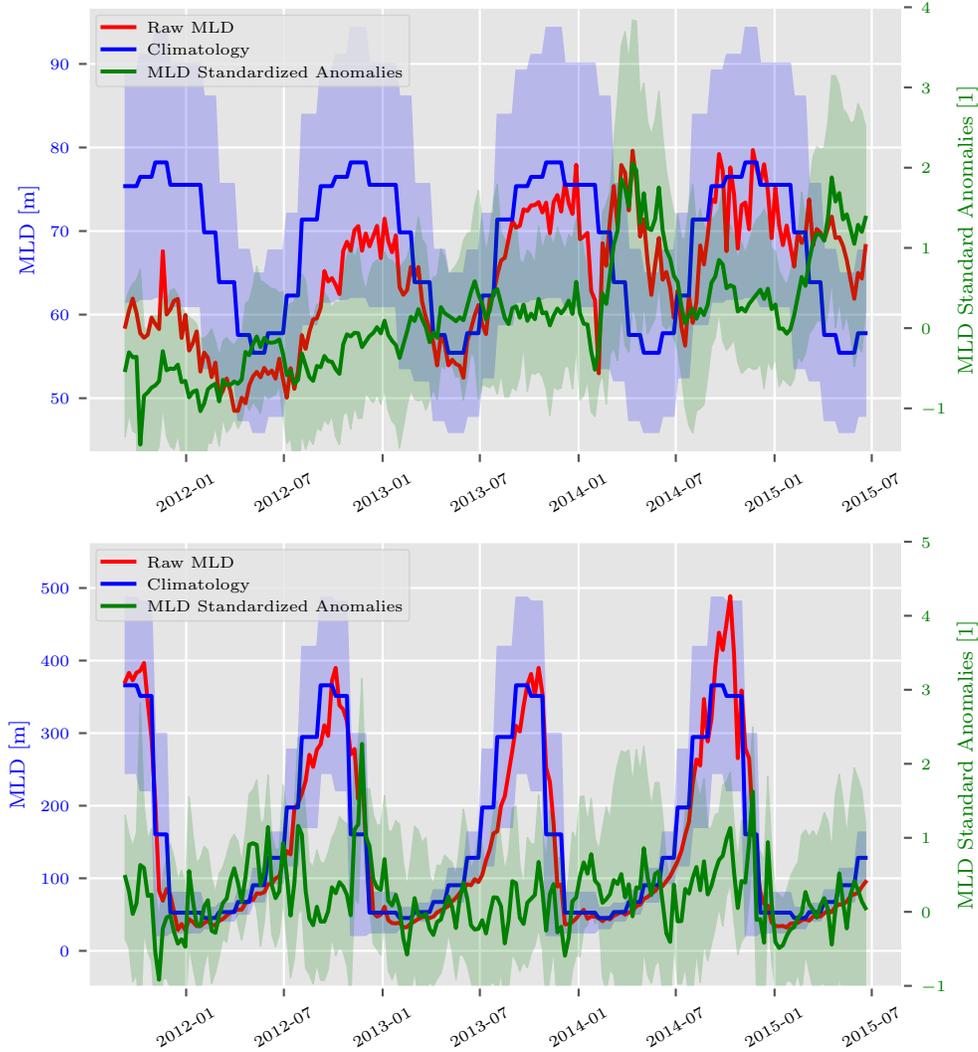


Figure 1. Several time series of the average MLD in each region at weekly resolution in the equatorial Pacific (top) and southern Indian Ocean (bottom), including the ensemble average of the MLD profiles over the domain (red), the ensemble average of the corresponding standardized MLD anomalies (green), and the area-average of the gridded monthly MLD climatology (blue). The blue shading represents the area-average of the gridded monthly standard deviations, and the green shading represents the ensemble standard deviation of the profile-wise standard anomalies. (Top) equatorial Pacific region (120W, 10S) - (150W, 10N). (Bottom) southern Indian Ocean (45S, 60E) - (35S, 120E).

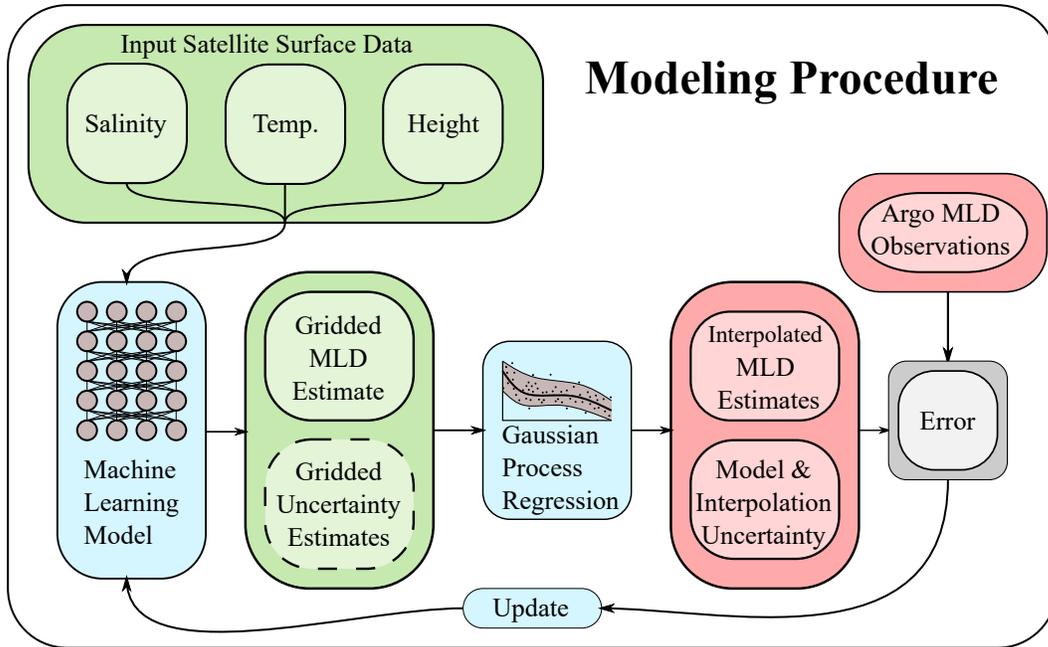


Figure 2. A schematic of the modeling procedure. Satellite sea surface data is fed into the machine learning model to produce a gridded MLD estimate (with some form of an uncertainty estimate if the machine learning model is probabilistic). To compare with the observations and optimize parameters, these gridded estimates are fed into a Gaussian process regression model (with its own hyper-parameters that are optimized) to produce MLD estimates interpolated to the locations where the Argo observations exist. These interpolated estimates are automatically associated with uncertainty estimates that derived from either just the Gaussian process interpolation uncertainty (if the model is deterministic) or a combination of the Gaussian process uncertainty with ML model uncertainty (if the ML model has uncertainty estimates). The interpolated estimates are then compared with the observations to estimate various errors.

213 and date. In addition, for each profile, we divide by the bin’s corresponding monthly standard
 214 standard deviations to create standardized anomalies. Fig. 1 shows the time series of the raw
 215 MLD data, including the ensemble average of the individual profiles in each region, the
 216 ensemble average of the standardized anomalies at each profile, and the area-average of
 217 the gridded climatology, in two spatial regions under study (120W, 10S) - (150W, 10N)
 218 and (45S, 60E) - (35S, 120E). The character of the anomalies and standardized anomalies
 219 are not dissimilar, but the standardized anomalies have a more appropriate scale for
 220 machine learning purposes (see the Acknowledgements for data availability).

221 2.5 Evaluation Regions

222 In order to evaluate the behavior of the machine learning models in two different
 223 oceanic regimes, we choose to investigate two geographic regions with very different MLD
 224 variability on timescales from subseasonal to interannual but significant subseasonal MLD
 225 variability to learn in both cases. First, we choose the equatorial Pacific Ocean (10°S
 226 - 10°N and 150°W - 120°W), which has modest subseasonal MLD standard deviations
 227 (~ 15 m), a small climatological annual cycle (~ 20 m), and substantial interannual
 228 variability (see Fig. 1 and (D. B. Whitt et al., 2019)). Second, we choose to study the
 229 southern Indian Ocean (45°S - 35°S and 60°E - 120°E), which features larger subsea-

230 sonal standard deviations (~ 50 m), a much larger climatological annual cycle (~ 300
231 m), but relatively weak interannual variability.

232 Hence, both regions contain substantial subseasonal MLD variability to learn, but
233 the absolute magnitudes of the subseasonal variability as well as the relative magnitudes
234 of subseasonal, seasonal, and inter-annual variability differ dramatically.

235 In order to test our framework for estimating MLD using sea surface information
236 we perform the following experiment on each region of interest. On the 150 (out of 200
237 total) weeks of training data, we apply the training procedure summarized in Fig. 2 and
238 described in more detail in section 3 (see the Acknowledgements for a link to the soft-
239 ware).

240 On the remaining 50 weeks of testing and validation data the model predicts a dense
241 grid of MLD estimates based solely on the sea surface information as input. From this
242 dense grid, we interpolate the estimates onto the locations where in-situ Argo profile ob-
243 servations of the MLD exist and compute error statistics between the interpolated es-
244 timates and the observations. The interpolation is done using a Gaussian process (see
245 section 3.1) regardless of the machine learning method. We denote this testing proce-
246 dure as measuring the out-of-sample performance of the method.

247 3 Methods

248 We consider a simple but general model for the relationship between the surface
249 variables, salinity (S), temperature (T), and height (H), and mixed layer depth model
250 output (d),

$$251 \quad d = f(S, T, H; \theta) + \sigma, \quad \sigma \sim \mathcal{N}(0, \Sigma). \quad (1)$$

252 where θ refers to the collection of function parameters. The surface variables ex-
253 ist on a pre-specified grid, \mathbf{x} , of total size M and the function f may generally couple
254 surface variables from across this grid to produce d at a particular grid point. The dif-
255 ference between the mixed layer and the output of f , σ , is assumed to be a normally dis-
256 tributed random variable according to the covariance Σ that expresses the spatial un-
257 certainties in this functional relationship. The exact structures and parameterizations
258 of f that we use in this paper are described in section 3.2 while the methods we use to
259 specify Σ are presented in section 3.3.

260 Both the functional relationship f and the covariance matrix Σ are data-driven (i.e.,
261 agnostic to the underlying physics) and informed via observations d_o that exist at ar-
262bitrary (ungridded) locations, \mathbf{x}_o where freely-drifting Argo floats collect a profile. In
263 order to couple the gridded variables with the ungridded observations, we define the re-
264 lationship between our model and the observations to be a Gaussian process,

$$265 \quad d_o = Ld + \nu, \quad \nu \sim \mathcal{N}(0, V), \quad (2)$$

266 which will be further defined in section 3.1. Importantly, L and V , the spatial pro-
267 jection and covariance matrices, are independent of the observation values and only de-
268 pend on the observation locations, model grid locations, and model uncertainties. The
269 Gaussian process relationship, in our study, is entirely a spatial relationship that accounts
270 for spatial covariance between observations of the MLD. This implicitly means, however,
271 that L and V change depending on the particular week the data is from, but only be-
272 cause the particular locations \mathbf{x}_o where estimation and validation occurs vary from week
273 to week.

274 A further consequence of the chosen relation between the observations and model
275 (2) is that it defines the objective function, i.e. the conditional likelihood probability dis-

276 tribution, that will be maximized to fit the parameters of the nonlinear function f :

$$277 \quad \ln p(d_o|d) = -\frac{1}{2}(d_o - Ld)^T V^{-1}(d_o - Ld) - \frac{1}{2} \ln |V| - \frac{M}{2} \ln 2\pi. \quad (3)$$

278 Details of this optimization procedure are given in section 3.2. Here, it is implic-
279 itly understood that d , and hence $p(d_o|d)$, is a function of the input variables S, T, H ,
280 the architecture of the function f , and the parameters of f, θ .

281 The Gaussian assumptions made in Eq. 1 is primarily for notational convenience.
282 The model definition (Eq. 1) can easily be modified to include non-Gaussian noise by
283 including a stochastic component in f , $f(S, T, H; \theta, \sigma)$. This type of noise component
284 is important if we expect the noise to be a nonlinear function of the surface variables.
285 To account for this possibility, two of the probabilistic machine learning methods that
286 we test in this paper, Dropout and Variational Auto-Encoders (see section 3.3) are for-
287 mally of this type and require sampling to determine the covariance for use in the Gaus-
288 sian process. The Gaussian assumption made in (Eq. 2) is a reflection of the belief that
289 the interpolating operator between the gridded locations and Argo locations is appro-
290 priately approximated by a linear function. We believe that this is not overly restrictive
291 since most optimal interpolation techniques make similar assumptions.

292 3.1 Gaussian Process Regression

293 Gaussian Process Regression is closely related to the somewhat more general Op-
294 timal Interpolation and Kriging frameworks. For a more detailed history and exposition,
295 see Cressie (1993). A Gaussian process is any collection of random variables for which
296 any finite number have a joint Gaussian distribution and, as a result, is completely de-
297 termined by a mean and covariance function (Rasmussen & Williams, 2006). Given a
298 set of (2-dimensional) observation locations $\mathbf{x} = (x_1, \dots, x_M)^T$, we define the mean func-
299 tion $m(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$ of the process $d(\mathbf{x})$ as

$$300 \quad m(\mathbf{x}) = \text{E}[d(\mathbf{x})] \quad (4)$$

$$301 \quad k(\mathbf{x}, \mathbf{x}') = \text{E}[(m(\mathbf{x}) - d(\mathbf{x}))(m(\mathbf{x}') - d(\mathbf{x}'))] \quad (5)$$

302 Typically the mean function is set to zero and covariance function is parameter-
303 ized according to some kernel function. Various kernel functions impart different types
304 of regularity (differentiability): the exponential kernel leads to non-differentiable out-
305 puts, the Matern Class of kernels have a regularity parameter, and the squared expo-
306 nential kernel leads to smooth outputs. In our study, the squared exponential kernel,

$$307 \quad k(\mathbf{x}, \mathbf{x}') = \alpha e^{-\frac{1}{2\ell} \|\mathbf{x} - \mathbf{x}'\|^2} + \beta \quad (6)$$

308 where α and ℓ are hyperparameters that control the amplitude and length-scale of the
309 corresponding covariance structure, was chosen because of its marginally better perfor-
310 mance and efficiency compared to Matern class kernels. We train our Gaussian process
311 hyperparameters by optimizing according to the Gaussian process prior probability dis-
312 tribution over the training observation points \mathbf{x} ,

$$313 \quad \ln p(\alpha, \ell, \beta|d) = -\frac{1}{2} d^T K(\mathbf{x}, \mathbf{x})^{-1} d - \frac{1}{2} \ln |K(\mathbf{x}, \mathbf{x})| - \frac{M}{2} \ln 2\pi, \quad (7)$$

314 where the covariance matrix has entries $K_{i,j}(\mathbf{x}, \mathbf{x}) = k(x_i, x_j)$. To regularize the op-
315 timization process and ensure positivity of α, ℓ , and β , priors are occasionally placed on
316 the hyperparameters in a Bayesian fashion. In our study, this type of implementation
317 had minimal impact on the optimized values. In circumstances where either computa-
318 tional considerations are not a concern or available training data is limited, it is also pos-
319 sible to optimize the hyperparameters by cross-validating and minimizing the conditional

320 likelihood distribution, for details see Rasmussen and Williams (2006). The variance hy-
 321 perparameter β can, in general, be made anisotropic at the expense of increasing the total
 322 number of hyperparameters, but we do not consider such options in this study.

323 During the training of the neural network, i.e. while optimizing the parameters in
 324 f via Eq. 3 using backpropagation on training data from a given week, the Gaussian pro-
 325 cess hyperparameters must be re-optimized according to Eq. 7 because the Gaussian pro-
 326 cess parameterization depends on the Argo profile locations (and model covariance Σ ,
 327 if available) which generally vary from one training week to the next.

328 Once the Gaussian process has been optimized using function values (\mathbf{x}, d) , we can
 329 perform inference at the Argo spatial locations \mathbf{x}_o to obtain estimates of d_o . The infer-
 330 ence procedure follows Eq. 2 with L and V given by the equations

$$331 \quad L = k(\mathbf{x}_o, \mathbf{x}) (k(\mathbf{x}, \mathbf{x}) + \Sigma)^{-1} \quad (8)$$

$$332 \quad V = k(\mathbf{x}_o, \mathbf{x}_o) - k(\mathbf{x}_o, \mathbf{x}) (k(\mathbf{x}, \mathbf{x}) + \Sigma)^{-1} k(\mathbf{x}, \mathbf{x}_o). \quad (9)$$

333 Thus, the trained kernel function is independent of time and depends only on distance
 334 $\|\mathbf{x} - \mathbf{x}'\|$ not location \mathbf{x} or time, but L and V depend on location and time because Σ
 335 depends on location \mathbf{x} and the particular points chosen for estimation \mathbf{x}_o (e.g., the Argo
 336 profiles locations) vary with time.

337 3.2 Machine Learning

338 The main objective of this paper is to learn a relationship between the sea surface
 339 variables (salinity, temperature, height) and mixed layer depth. Without an a priori physics-
 340 based model, one must choose a reasonably parameterized model to approximate this
 341 relationship. Traditionally this relationship is represented via some linear or simple non-
 342 linear parameterization where one hopes that the true relationship lies in, or is not too
 343 far from, the output space of the model. For example, a basic linear model that we test
 344 in this paper is of the form,

$$345 \quad d_\ell = \begin{bmatrix} c_1(\mathbf{x}) \\ c_2(\mathbf{x}) \\ c_3(\mathbf{x}) \end{bmatrix} \cdot \begin{bmatrix} S \\ T \\ H \end{bmatrix} + b + \sigma, \quad \sigma \sim N(0, \Sigma) \quad (10)$$

346 Such models, however, are typically not expressive enough to represent arbitrary
 347 relationships. The revolution of machine learning, and, in particular, deep learning, has
 348 been born out of the need to express arbitrary functional relationships amid a dearth
 349 of observational data. While there exists several popular machine learning architectures,
 350 we base our paper around modifications of the quintessential deep learning model, the
 351 feedforward neural network (FNN) (Goodfellow et al., 2016). FNNs are represented by
 352 composing together many different functions in series to form a chain,

$$353 \quad f(x) = f^{(n)}(f^{(n-1)}(\dots f^{(1)}(x) \dots)), \quad (11)$$

$$354 \quad f^{(i)}(x) = a(x^T W_i + b_i), \quad (12)$$

355 where W_i is a matrix of weights, b_i is a bias term, and $a(\cdot)$ is what is referred to
 356 as an ‘activation function’, that applies a simple non-linearity element-wise to the affine
 357 transformation of the input, x . Common examples of activation functions include the
 358 sigmoid, softplus, and rectified linear functions. Based on the experiments in Gal (2016),
 359 we implement the rectified linear unit as the activation function in all of our neural net-
 360 work layers, although it is possible that, among all of the available activation functions,
 361 another function would result in superior performance. We will denote the collection of
 362 neural network parameters as $\theta = \{W_1, \dots, W_n, b_1, \dots, b_n\}$.

The training of a neural network entails obtaining an estimate of the parameters, $\hat{\theta}$, by approximately solving the optimization problem,

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} \ln p(d_o|d) \\ &= \arg \min_{\theta} \left\{ g(\theta) - \sum_{j=1}^{n_{\text{train}}} \ln p_j(d_o|d) \right\} \end{aligned} \quad (13)$$

where $g(\theta)$ is a regularization function that is applied to both constrain the possible parameter values and stabilize the optimization procedure. As written, $p_j(d_o|d)$ refers to the joint probability distribution between the j^{th} input and output data. The optimization procedure includes all training data but, in practice, subsetting is common (as in batch gradient descent (Ruder, 2016)). We only seek an approximate solution to Eq. 13 for two reasons: first, the optimization problem is highly-non trivial, non-convex, and high-dimensional with many local minima and obtaining a global minimum is infeasible; second, the ultimate goal is for the parameters to lead to a function f that generalizes well to data not in the training set and over-training might ultimately hinder this goal (Caruana et al., 2001). The problem of over-fitting and poor generalization is one of the largest obstacles to good machine learning performance, particularly in applications where prediction involves extrapolation beyond whatever data was in the training set. All of the neural networks implemented for this paper are done using the TensorFlow and TensorFlow Probability frameworks (Abadi et al., 2016; Dillon et al., 2017).

Because our study is limited to only 150 training weeks, we implement a non-standard training strategy to help reduce overfitting. For each epoch (a single run through the entire training data) we divide the 150 training weeks randomly into 6 batches of 25 weeks. The first batch is held out and the current loss on that batch is saved. For each subsequent batch, the loss for that batch is used to update the model parameters. To update the parameters, we use the Adam optimizer with initial learning parameter set to $1e-3$ (Kingma & Ba, 2015). With the updated model parameters, we calculate a new loss on the first, held-out batch. If that new loss is less than the saved loss, then the updated parameters are accepted and the new loss is saved. If the new loss is larger than the saved loss then the parameters are only accepted with

$$\text{probability of acceptance} = \exp(\text{saved loss} - \text{final loss}).$$

This training strategy reduces the amount of overfitting because it forces updates to be generalizable to the held out batch, which acts as a 'testing batch'.

FNNs with enough hidden layers have been proven to serve as a universal approximator (Hornik et al., 1989; Cybenko, 1989; Leshno et al., 1993). This means that, at least theoretically, there exists a FNN that can represent whatever functional relationship exists between the sea surface variables and MLD. Unfortunately, there is no guaranteed way to find this optimal relationship. While the optimization problem (Eq. 13) has a natural inherited probabilistic framework, even an exact solution has no guarantee of agreeing with the 'true' relationship. The construction of these optimization frameworks and the regularization functions is often done by trial and error since there is, as of yet, no clear casual relationship between tuning the architecture settings and the resulting uncertainty estimate - even if the model can be viewed through a (Bayesian) probabilistic framework.

Finally, since the (approximate) solution to Eq. 13 is not accompanied with natural uncertainty estimates for the parameters, it can be difficult to obtain calibrated probabilistic estimates of \hat{d} . To truly obtain samples from the posterior $p(d|d_o, S, T, H, \theta)$, we would need to incorporate any and all uncertainties that exist in the inputs, observations, model parameters, and model framework and be able to sample from them effectively. Due to the high-dimensionality of the problem, this is computationally infeasible and therefore we must rely on adequate approximations. In the next section, we outline the approximations that we test in this manuscript.

3.3 Probabilistic Machine Learning Models

The simplest technique to introduce uncertainty estimates into a neural network is to implement Dropout (Hinton et al., 2012; Srivastava et al., 2014). Acting as a layer of the network, Dropout randomly sets inputs to zero at a particular rate and scales the rest of the inputs by $1/(1 - \text{rate})$. Mathematically,

$$f^{(i)}(x) = \frac{1}{1-p} M \odot a(x^T W_i + b_i), \quad M_j \sim \text{Bernoulli}(p), \quad (14)$$

where \odot means element-wise multiplication. Each run of the model then has a different combination of weights that are set to zero. While originally this technique was used to reduce overfitting, it can also be viewed through a Bayesian probabilistic lens (Maeda, 2014). Running the model multiple times creates an ensemble that can be used to calculate moments of the output distribution, and, in particular, Σ and μ . It has been shown that the expected distribution from a neural network utilizing Dropout forms a Gaussian mixture distribution (Gal & Ghahramani, 2016). Therefore, there is some reason to believe that the regularity of the data distribution dictates how useful Dropout can be in uncertainty quantification.

The next simplest probabilistic technique, what we call the Variational Artificial Neural Network (VANN), also known as a heteroscedastic network, is to parameterize the output of the neural network according to some distribution. For a Gaussian distribution, for example, the output of f is a stacked vector of the mean and covariance estimates,

$$f(S, T, H; \theta) = [\mu; \text{vec}(\Sigma)], \quad (15)$$

where $\text{vec}(\Sigma)$ is the flattened covariance matrix, such that $d \sim N(\mu, \Sigma)$. This technique is relatively easy to implement with care needed to ensure that constraints on the parameters are enforced. Typically, a Bayesian framework would then impose prior probability distributions onto μ and Σ . In particular, in addition to the Gaussian likelihood, it is common to impose a Gamma or LKJ - uniform over the space of covariance matrices - prior on the covariance to prevent unnecessary shrinkage. In a feedforward neural network, this parameterization increases the number of outputs and hence the overall total number of parameters. If the number of grid points of $d(\mathbf{x})$ is M then a full covariance matrix would require $M(M + 1)/2$ parameters and the corresponding number of parameters required in the neural network makes it computationally prohibitive as k grows large. To limit the computational cost, we make a diagonal assumption about the covariance to reduce the number of parameters at the expense of losing covariance information between MLD values at different grid points. Parameterization of the data distribution is not always possible if a good approximation or transformation to an appropriate probability distribution is not known and the effectiveness of this technique is reflection of the quality of that assumption.

The final method that we test is the variational auto-encoder (VAE) (Kingma & Welling, 2014). A typical VAE consists of two dense networks: an encoder that projects the inputs into a lower-dimensional latent space, parameterized by a probability distribution, and a decoder that inverts this projection and produces the original input. The loss between the decoder's output and the original system drives the learning process. A VAE supposes a prior distribution over the latent variable z , $p(z)$, that, along with the decoder network that induces a conditional likelihood distribution $p(S, T, H|z; \theta)$, forms a posterior distribution,

$$p(z|S, T, H; \theta) \propto p(z)p(S, T, H|z; \theta)$$

This posterior distribution is typically intractable and thus replaced by a variational approximation $q(z|S, T, H; \theta)$. This approximation includes a parameterization of the prior and likelihood distributions, typically Gaussian distributions with parameters that are

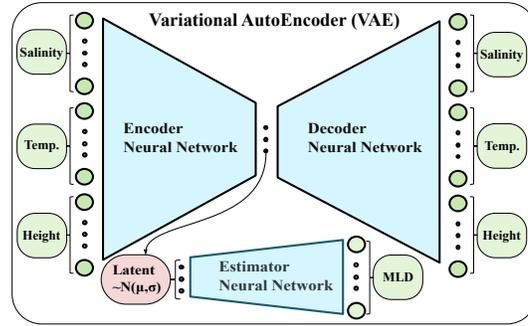


Figure 3. A schematic of the modified VAE. Training is informed by the decoder and estimator networks losses. For a full description of the training procedure for a typical VAE, see Kingma and Welling (2014).

462 learned in the encoder network. In our design we also use a Gaussian distribution in the
 463 latent space, and, as demonstrated in Figure 3, we couple this network with a third dense
 464 network, which we call the estimator, that transforms the latent space into an estimate
 465 of the MLD associated with the surface salinity, temperature, and sea height anomaly
 466 encoder inputs.

467 While the prior and likelihood distributions in a VAE are specified as Gaussian,
 468 the distribution of the output of the estimator network, that is, the MLD outputs, is not
 469 parameterized. While the difference between the MLD estimates and the MLD obser-
 470 vations is modelled as a Gaussian process regardless of neural network architecture, the
 471 possible benefit of our chosen VAE approach is that it can produce theoretically arbi-
 472 trary probability distribution $p(d|S, T, H; \theta)$. Another theoretical benefit to this approach
 473 is that, since the neural network can learn an efficient lower-dimensional representation
 474 of the inputs that capture dominant patterns, the estimator might be better able to gen-
 475 eralize and less sensitive to small perturbations and noise in the inputs.

476 We summarize the ways in which the MLD uncertainty, represented as Σ , is esti-
 477 mated. For the non probabilistic methods (linear model, artificial neural network), there
 478 is no associated Σ . For the Variational Artificial Neural Network (VANN), Σ is a direct
 479 output of the neural network and the weights that produce this Σ are trained as in Eq.
 480 13. For the Dropout network, each output of the network is a draw from a random dis-
 481 tribution. Σ is the sample covariance matrix of 100 random samples from this distribu-
 482 tion. Similarly, for the variational auto-encoder (VAE), Σ is the sample covariance ma-
 483 trix from 100 random outputs of the VAE network.

484 4 Experimental Results

485 We test 6 different methods on each experiment, five of which we consider as part
 486 of the machine learning framework: the linear model (Eq. 10), the feedforward artificial
 487 neural network (Eq. 11), feedforward neural network with parameterized distributional
 488 output (Eq. 15) feedforward neural network with Dropout (Eq. 14), and a variational
 489 auto-encoder. We collectively shorthand these to be 'Linear', 'ANN', 'VANN', 'Dropout',
 490 and 'VAE'. While the models presented in this study are based on the basic feedforward
 491 neural network architecture, we also tested (with poor performance) convolutional neu-
 492 ral networks with a multitude of architectures and hyperparameters. Finally, in order
 493 to compare these methods to a traditional interpolation only approach, we implement
 494 an Ordinary Kriging scheme, which we call 'OI' for optimal interpolation, with a (spa-
 495 tial) spherical kernel chosen via cross-validation and parameters optimized via maximum

596 likelihood. The OI approach only uses the in-situ MLD standard anomaly observations,
 597 with no sea surface information, to make gridded estimates. Therefore, even during the
 598 out-of-sample prediction experiments, the OI’s error statistics for a given week are cal-
 599 culated using only that week’s data. In particular, we use a cross-validation approach
 600 using a 75-25% train-test split to estimate these error statistics.

601 We use 3 metrics in our testing: root mean squared error (RMSE), Pearson cor-
 602 relation coefficient, and probabilistic calibration. These metrics are applied to modeled
 603 standardized MLD anomalies at the validating Argo profile locations (see section 2 for
 604 details). We use the typical definition of root mean squared error,

$$605 \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n |(d_o)_i - L(d)_i|^2}. \quad (16)$$

606 RMSE is a convenient metric in that it captures the mean prediction error, but it
 607 doesn’t necessarily tell us much about the relationship between the predictions and ob-
 608 servations and it also fails to capture meaningful information about the uncertainty of
 609 the predictions. To compensate for the first deficiency, we rely on the Pearson correla-
 610 tion coefficient (correlation) to provide insight into the existence of (linear) relationships
 611 between predictions and the Argo MLD data. For reference, correlation is defined as

$$612 \quad \text{Correlation} = \frac{\sum_{i=1}^n \left(L(d)_i - \overline{L(d)} \right) \left((d_o)_i - \overline{d_o} \right)}{\sqrt{\sum_{i=1}^n \left(L(d)_i - \overline{L(d)} \right)^2} \sqrt{\sum_{i=1}^n \left((d_o)_i - \overline{d_o} \right)^2}} \quad (17)$$

613 Common metrics that capture probabilistic calibration include skill scores such as the
 614 Brier score or the Kolmogorov–Smirnov statistic. Here, for simplicity, convenience, and
 615 data-limitation reasons, we use the following measure for probabilistic calibration,

$$616 \quad \text{Calibration} = \frac{1}{n} \sum_{i=1}^n \mathbb{1} \left[|(d_o)_i - L(d)_i| < \sqrt{V_{ii}} \right], \quad (18)$$

617 where V_{ii} is the i th diagonal entry of the covariance matrix of the Gaussian pro-
 618 cess regressor (Eq. 8) and $\mathbb{1}$ is 1 if the argument is true and 0 otherwise. Calibration is
 619 then a number between 0 and 1. It is important to remember that V also includes the
 620 covariance estimate from the probabilistic machine learning models, Σ . For non-probabilistic
 621 machine learning model, V does not include any model uncertainty beyond the learned
 622 hyperparameter β in Eq. 6. For a Gaussian statistic, the Calibration is theoretically \approx
 623 0.68, the optimal score for this metric. If a model scores lower than that theoretical thresh-
 624 old, it is underestimating the amount of uncertainty in the data. Conversely, a higher
 625 Calibration than the theoretical threshold represents an overestimation of the uncertainty.

626 We aim to give an overview of the main results from our studies. We focus on the
 627 aforementioned metrics as we compare model performance overall, and broken down by
 628 groups representing different levels of standard deviation in the observations. These met-
 629 rics indicate 3 conclusions: 1) Model performance is superior in the equatorial Pacific
 630 Ocean than the southern Indian Ocean, 2) the probabilistic machine learning methods
 631 outperform traditional OI, particularly in terms of correlation and calibration, and, there-
 632 fore, 3) the relative performance of machine learning algorithms indicate that surface vari-
 633 ables can provide meaningful information about the mixed layer depth and produce es-
 634 timates that are as good or better than OI methods that directly use MLD data. Finally,
 635 we visually compare the model outputs in two case studies that represent the best and
 636 worst model performance. To provide context and further applications, we also include

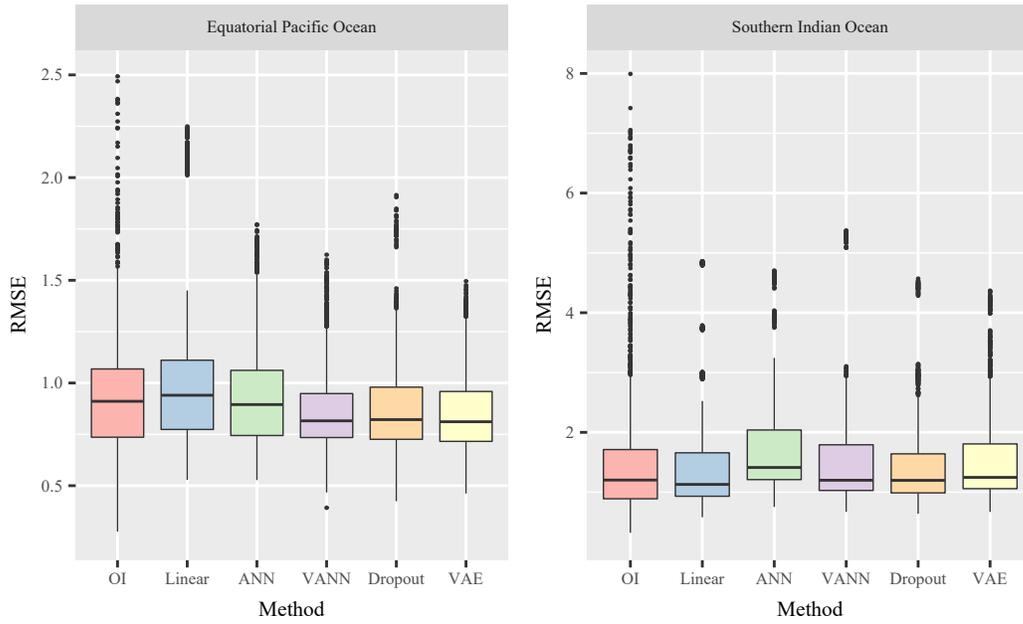


Figure 4. Root Mean Squared Errors (RMSE) on temporal out-of-sample prediction (in meters). Errors are calculated on 50 withheld validation weeks. Boxes capture 25-75% of the weekly errors with the middle line representing the median error. Dots are considered outliers - values which are $1.5\times$ lower/upper quantile. (Left) The equatorial Pacific region (120W, 10S) - (150W, 10N). (Right) The southern Indian Ocean region (45S, 60E) - (35S, 120E). Note the difference in scales between the two regions. OI errors are calculated using cross-validation within each week (see text for details).

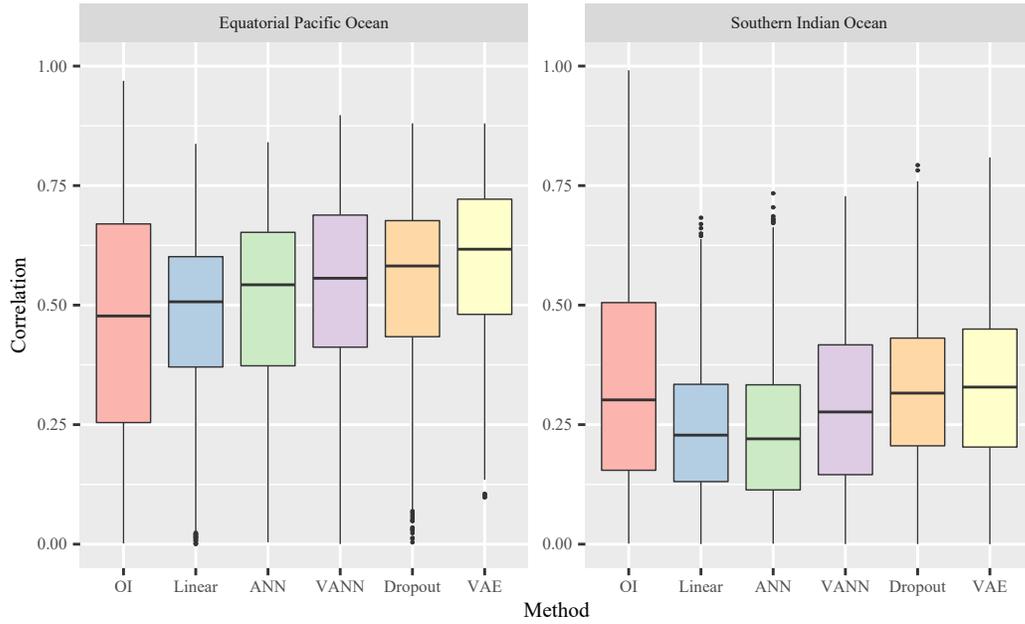


Figure 5. Correlation on temporal out-of-sample prediction (in meters) as in Fig. 6. Correlations are calculated on 50 withheld validation weeks. Boxes capture 25-75% of the weekly correlation with the middle line representing the median correlation. Dots are considered outliers - values which are $1.5\times$ lower/upper quantile. (Left) The equatorial Pacific region (120W, 10S) - (150W, 10N). (Right) The southern Indian Ocean region (45S, 60E) - (35S, 120E). OI values are calculated using cross-validation within each week (see text for details).

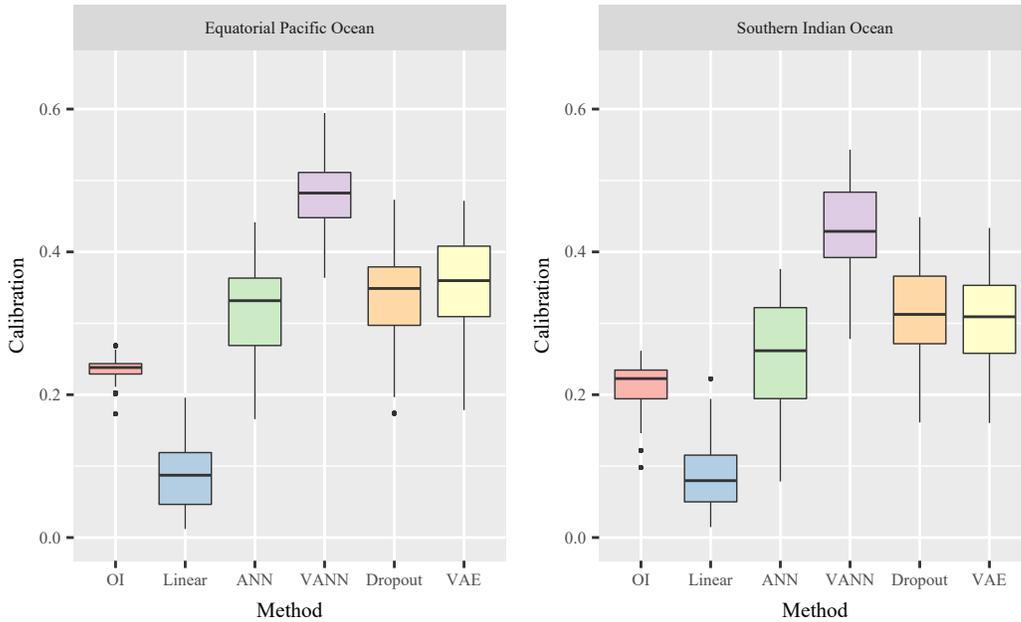


Figure 6. Measure of probabilistic calibration on temporal out-of-sample prediction as in Fig. 6. Calibrations are calculated on 50 withheld validation weeks. For each week, we find the percent observations that fall within 1 standard deviation of forecast ensembles. For a Gaussian distribution, this probability should be approximately 0.68, with greater relative values representing under-confident and lesser relative values representing overconfident predictions. OI calibrations are calculated using cross-validation within each week. (Left) The equatorial Pacific region (120W, 10S) - (150W, 10N). (Right) The southern Indian Ocean region (45S, 60E) - (35S, 120E).

537 model outputs from the HYCOM + NCODA Global 1/12° Analysis (Fox et al., 2002;
 538 Cummings, 2006; Cummings & Smedstad, 2013) for a visual comparison with our purely
 539 data-driven approaches.

540 Considering first the RMSE, model performance is superior in the equatorial Pa-
 541 cific Ocean compared to the southern Indian Ocean, and the various models differ only
 542 modestly within each region. Fig. 4 (note the difference in scales of the vertical axis) shows
 543 the RMSE results over the two regions. The machine learning methods seemingly per-
 544 form well against OI, particularly in the equatorial Pacific as the Dropout and VAE meth-
 545 ods have the lowest median RMSE and 25% - 75% range. In the southern Indian Ocean,
 546 the Linear method performs well, initially suggesting that the mean dynamics can be
 547 well approximated by a linear combination of the surface variables. The number and range
 548 of OI outliers, in comparison to machine learning approaches, demonstrates that the ma-
 549 chine learning approaches offer more stable predictions.

550 The correlation analysis underscores and further confirms the result (derived from
 551 RMSE above) that the overall model performance is better in the eastern equatorial Pa-
 552 cific Ocean compared to the southern Indian Ocean (Fig. 5). However, the correlations
 553 also reveal more substantial differences between the models in each region. In the equa-
 554 torial Pacific, it is clear that the machine learning methods perform better than tradi-
 555 tional OI, with the VAE performing the best. In the southern Indian Ocean, however,
 556 there is little separating the performance between OI and probabilistic machine learn-
 557 ing methods, although the VAE is marginally the best performing model in this region
 558 as well. A key difference between the RMSE results in Fig. 4 and the correlations in Fig.
 559 5 is that the linear method, while having a small predictive RMSE, has poor correlation
 560 with the observations. From other testing, we believe that the linear model has both small
 561 RMSE and correlation because the outputs of the linear method are generally smaller
 562 values.

563 The calibration results in Fig. 6 demonstrate that the probabilistic machine learn-
 564 ing approaches using surface data are significantly better at estimating the posterior un-
 565 certainty than OI and MLD data alone. Furthermore, model performance is again su-
 566 perior (albeit modestly so) in equatorial Pacific Ocean compared to the southern Indian
 567 Ocean. The linear model performs very poorly in comparison to the other machine learn-
 568 ing methods. The traditional OI approach also has poorer performance compared to the
 569 machine learning models. In addition, all probabilistic techniques appear to perform slightly
 570 better than the non-probabilistic ANN (in terms of both calibration and RMSE). How-
 571 ever, the smallness of the differences between ANN and the other ML models suggests
 572 that much of the uncertainty manifest in all the ML model calibrations is due to the Gaus-
 573 sian Process regression, since the ANN does not have inherent MLD uncertainty esti-
 574 mates. Among the three probabilistic machine learning models, VANN, Dropout, and
 575 VAE, the VANN has dramatically better calibration than the other two methods. This
 576 discrepancy shows that, in these particular case studies, explicitly parameterizing the
 577 noise better captures the underlying uncertainty than the sampling-based approaches.

578 The conclusion from the calibration metric are mirrored in Fig. 7, where the em-
 579 pirical cumulative distribution of the models is plotted against the distribution of the
 580 observations. The diagram represents the Lines closer to the optimal red line in that fig-
 581 ure represents better model calibration. It is clear from this plot that the VANN and VAE
 582 have superior performance in estimating the tails of the distribution when compared to
 583 other methods and the OI. It is true, however, that overall performance is lacking. The
 584 behavior of each line indicates that the tails of the model distribution are shorter than
 585 the observational distribution - another indication that extreme MLD values remain dif-
 586 ficult for the models to predict.

587 The difference between performance in VANN vs. Dropout and VAE could plau-
 588 sibly explained by suggesting that the posterior probability distribution of the MLD given

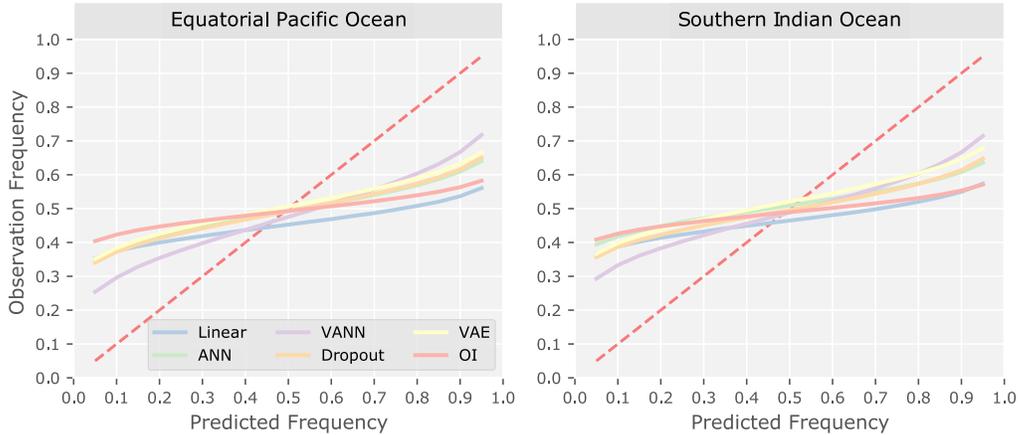


Figure 7. Probability plot comparing the empirical cumulative distributions of the model outputs against the data. The dotted red line would represent perfect agreement between models and observations. A value above and to the left of the red line indicates a part of the distribution that is over-represented, whereas a value below and to the right of the red line indicates a part of the distribution that is underrepresented.

589 satellite data is closely approximates a Gaussian distribution and hence well estimated
 590 by the VANN. Alternatively, the available data may not be sufficient to allow the sampling-
 591 based methods (Dropout and VAE) to learn the posterior distribution.

592 To reveal how the model performance depends on the MLD variability, we group
 593 the observed MLDs at the Argo profile locations by the (ensemble) standard deviation
 594 of all observed standardized MLD anomalies (defined in section 2.4) in the same week
 595 and region using K-Means clustering. We find that model performance generally degrades
 596 in terms of RMSE (Fig. 8) but improves in terms of correlation (Fig. 9) in weeks with
 597 higher standard deviations. But, model calibration (not shown) is relatively insensitive
 598 to the weekly variability of MLD anomalies. With regard to RMSEs in Fig. 8, we find
 599 that the increases in RMSE with standard deviation are fairly consistent across the mod-
 600 els, and the slope RMSE-over-standard-deviation is roughly 1 in both regions. In addi-
 601 tion, the probabilistic machine learning models have about equal or smaller RMSE than
 602 the OI at all levels of variance. Finally, it is notable that for the weeks with the largest
 603 observation standard deviations, the OI has particularly large RMSEs in the southern
 604 Indian Ocean, whereas the linear method has particularly large RMSEs in the equator-
 605 ial Pacific.

606 With regard to the correlations in Fig. 9, we find that the increasing standard devi-
 607 ation of the observations in the equatorial Pacific Ocean improves model performance
 608 to a much greater degree than in the southern Indian Ocean. Interestingly, the compar-
 609 isons between the models within each standard deviation cluster qualitatively mirror those
 610 of the whole dataset (c.f., Figs. 9 and 5): machine learning models generally produce higher
 611 correlation than OI, particularly in the equatorial Pacific Ocean. The only notable excep-
 612 tion is the bin with high standard deviation in the Southern Indian Ocean, where the
 613 VANN, Dropout and VAE models have notably higher correlation than the other meth-
 614 ods, while OI performs particularly poorly. Finally, the relatively high correlations at
 615 large standard deviation in the equatorial Pacific suggest, potentially, that the dynam-
 616 ics that cause large mixed layer depth anomalies also strongly couple with the surface
 617 variables in this region.

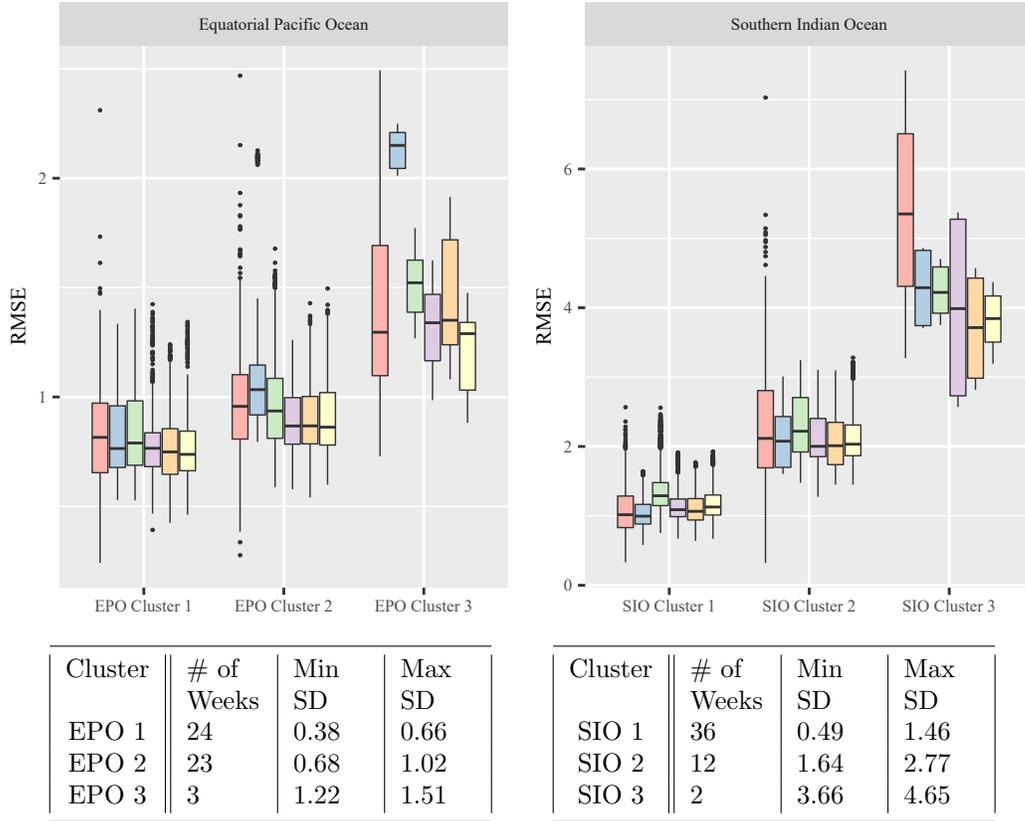


Figure 8. RMSEs divided by region and clustered by the standard deviations of the ensembles of MLD standard anomalies in a given week in (Left) the equatorial Pacific Ocean and (Right) the southern Indian Ocean. (Bottom) Table showing the number of weeks in each cluster, the minimum standard deviation in each cluster, and the maximum standard deviation in each cluster. (Top) The distribution of the RMSE for each method, corresponding to 40 samples from the posterior distribution for each week, separated by cluster. The boxplots are colored as in Fig. 4. Note the difference in scales between the two regions.

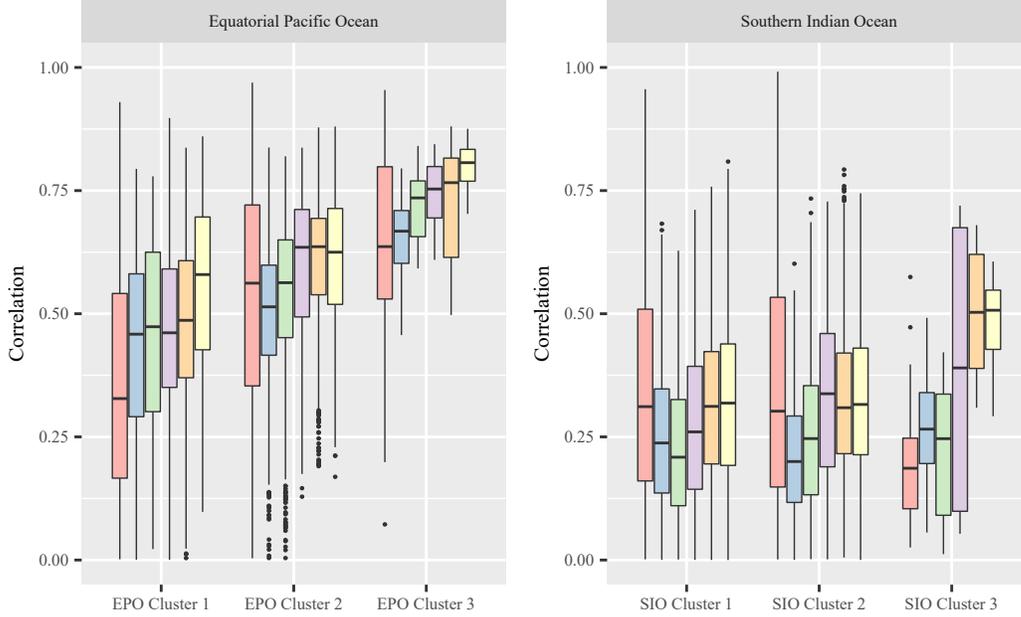


Figure 9. As in Fig. 8, but correlations instead of RMSE.

618 Taken together, the results indicate that, in the equatorial Pacific Ocean and to
 619 a lesser extent in the southern Indian Ocean, the surface information provides just as,
 620 if not more so, valuable information in estimating the MLD as the existing Argo obser-
 621 vations of the MLD.

622 To give a visual and spatial sense of the range of model estimates, we demonstrate
 623 two extreme ends of the prediction spectrum, the worst and best predictive weeks for
 624 our models. Full model output for all available weeks are available online (Foster et al.,
 625 2020) at <https://www.doi.org/10.5281/zenodo.4421752>. The week corresponding to
 626 the worst RMSE performance is the week of 11-23-2012 in the southern Indian Ocean.
 627 If you compare this with Fig. 1, this corresponds to a period of particularly large anomalies.
 628 The average RMSEs for this particular week corresponding to the OI and VAE mod-
 629 els are approximately 6.14 and 4.16. Similarly, the week of the relative best performance
 630 (now in terms of correlation coefficient) is 05-09-2014 in the equatorial Pacific Ocean.
 631 The corresponding average correlations (and RMSEs) for the OI and VAE methods are
 632 0.68 (1.09) and 0.83 (0.99). Figs. 10 and 11 show A. the data with overlaid sea level height
 633 contours, B. smooth gridded climatology, C. standard anomaly OI model output, D. VAE
 634 model output, E. reanalysis of VAE model output and observations, and F. HYCOM+NCODA
 635 Global 1/12° Analysis for these two weeks. MLD values are derived from the HYCOM+NCODA
 636 reanalysis by applying the MLD definition in D. B. Whitt et al. (2019) and averaging
 637 over the appropriate week (the raw southern ocean HYCOM data is 3-hourly and the
 638 equatorial Pacific data is daily). Each of the machine learning and OI model outputs are
 639 computed as MLD standard anomalies and are transformed back to MLD estimates for
 640 plotting. Because the output of the VAE model does not use observations at prediction
 641 time, we can perform our own reanalysis by finding the minimum of the associated pos-
 642 terior distribution,

$$\begin{aligned}
 \hat{d} &= \arg \min_d -\ln p(d|d_o, d_m), \\
 &= \arg \min_d (d - d_m)^T \Sigma^{-1} (d - d_m) + (Ld - d_o)^T V^{-1} (Ld - d_o).
 \end{aligned}
 \tag{19}$$

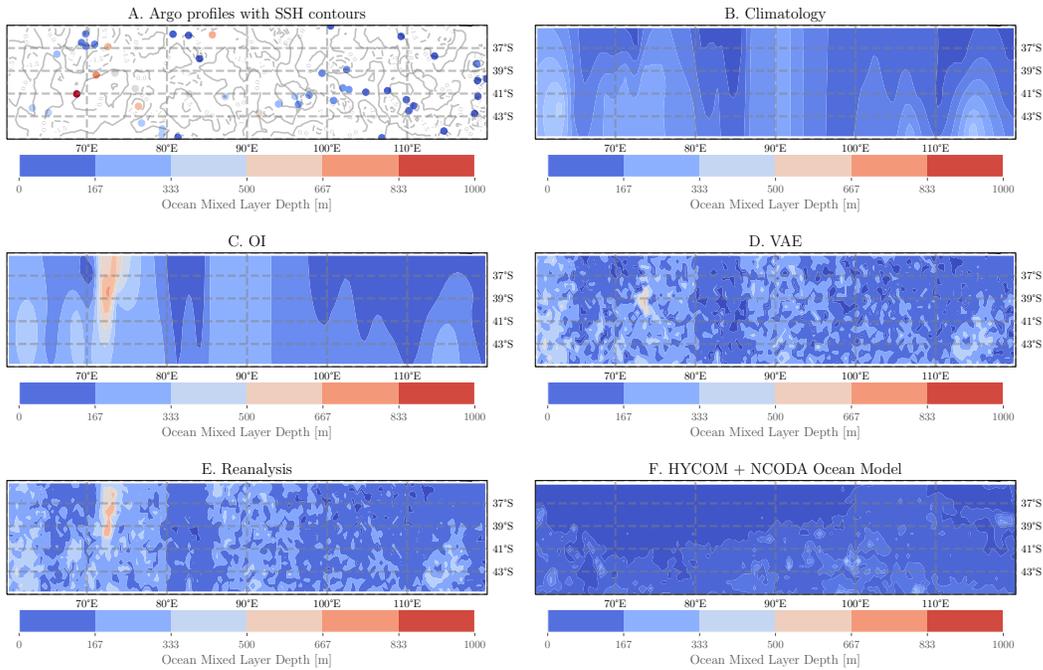


Figure 10. MLD estimates, estimated on standard anomalies with climatologies added back in, corresponding to the date of worst RMSE, achieved by VAE approach in the southern Indian Ocean, 11-23-2012. Methods from top left to bottom right: A. Argo float observations with sea level height contours of 0.5 meters are overlaid (blue is lower height), B. smooth gridded climatology, C. optimally interpolated standard anomalies with climatologies, D. VAE model with climatologies, E. Reanalysis of VAE and observations, and F. HYCOM+NCODA ocean model - see text for more details.

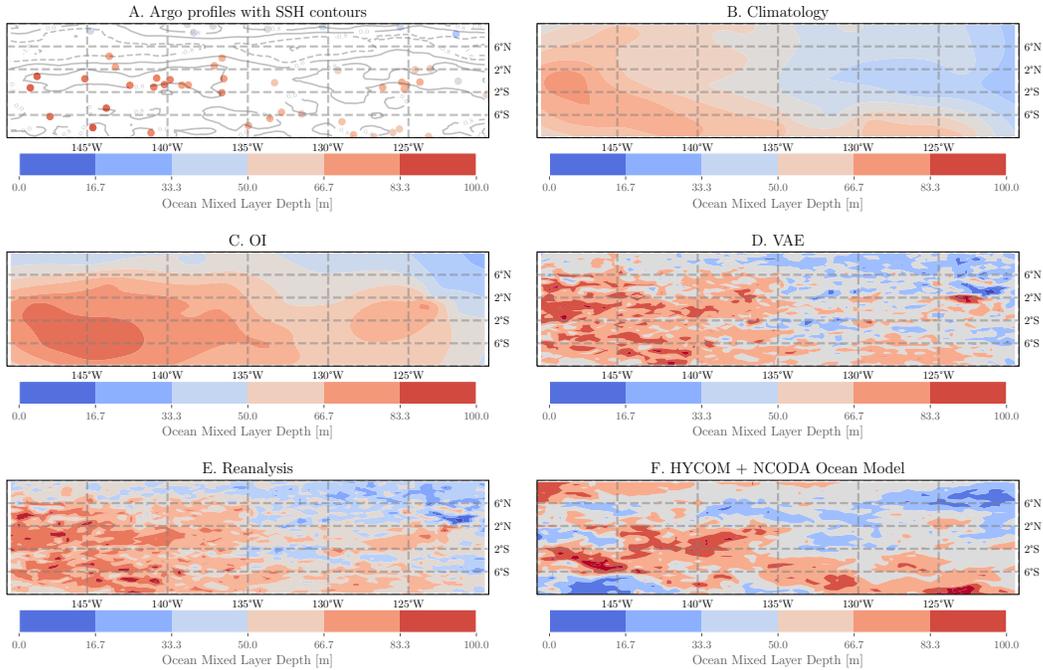


Figure 11. MLD estimates, estimated on standard anomalies with climatologies added back in, corresponding to the date of best average correlation, achieved by VAE approach in the equatorial Pacific Ocean, 05-9-2014. Methods from top left to bottom right as in Fig. 10

644 In Fig. 10, the week representing the collectively worst model performance, is an
 645 example of an extremely large MLD standard anomalies that can occur in late spring
 646 due to a delay in the springtime transition from deep winter to shallow summertime MLDs,
 647 as seen in Fig. 1. In this week, there is a narrow cluster of abnormally large MLD values
 648 that are visible in panels A, C, D and E. The OI model outputs are visually smooth,
 649 as a result of the spherical kernel used to do the interpolation, but underestimate the
 650 magnitude of the data. The VAE model output, as a result of being a function of the
 651 sea surface data, contains many small scale features that create a visually noisy gridded
 652 estimate. In addition, there are clusters of large anomalies where the data does not suggest
 653 any (near 115°E and 43°S for example). The reanalysis, as a result of being a variance-
 654 weighted average between the VAE and the observations, more closely resembles the OI
 655 estimate but still contains much more small scale variability. In the HYCOM + NCODA
 656 Global reanalysis, the model does not seem to capture the large MLD values that are
 657 seen in the Argo data, which might be due to the relative uncertainties in the HYCOM
 658 + NCODA Data Assimilation procedure. Direct comparisons between the VAE reanalysis
 659 and HYCOM+NCODA model should not be over-exaggerated because the differences
 660 in variance specification.

661 Similar to the worst case, the best case (achieved by the VAE model) occurs in a
 662 week of large standard anomalies in the equatorial Pacific (Fig. 1). As opposed to the
 663 worst case study, in this case study (Fig. 11) the climatology offers a lot of structure that
 664 is manifested in the MLD that week. The OI model output presents a very spatially coherent
 665 MLD estimate. The machine learning models, as a result of being functions of
 666 the sea surface inputs, have smaller scale features that modify the overall structure of
 667 the gridded MLD. The VAE model output, while having better performance in estimating
 668 the MLD standard anomalies than the OI at the observation locations, appears to
 669 have a greater stratified estimate. That is, the VAE model seems to overestimate the mag-

670 nitude of standard anomalies. The reanalysis of the VAE model output and observations
 671 retains a mixture of the smaller scale feature from the VAE model and the coherent struc-
 672 ture apparent in the OI output. The HYCOM + NCODA reanalysis closely captures the
 673 scale of the Argo MLD values, but the overall structure does not visually seem to match
 674 the observations. Again, the comparison with the HYCOM + NCODA reanalysis should
 675 be taken with appropriate qualification.

676 5 Conclusion and Discussion

677 The ocean mixed layer interacts with the atmosphere and deep ocean on a mul-
 678 titude of spatial and temporal scales. Heat exchange between these bodies has signifi-
 679 cant impact on subseasonal and interannual (aseasonal) timescales and can influence the
 680 behavior of dominant modes of variability (i.e. ENSO, MJO, tropical cyclones). Prolif-
 681 eration of Argo floats have dramatically increased the number of observations of the ocean
 682 over the preceding decades but are still too sparse to resolve fine spatio-temporal fea-
 683 tures of the MLD. Satellite data, however, is able to provide fine resolution gridded maps
 684 of sea surface variables, but cannot provide subsurface information.

685 The first goal of this work was to analyze the extent to which satellite data of sea
 686 surface variables can provide information useful for estimating the MLD. We built sev-
 687 eral machine learning models to learn such a relationship based on available data. We
 688 found that in terms of both root mean squared error, correlation, and probabilistic cali-
 689 bration, the machine learning model results suggest that the satellite data is equally if
 690 not more useful in estimating MLD values and uncertainties than MLD observations alone,
 691 given that sufficient MLD observations are available for out of sample training (Figs. 4
 692 & 6). The exact relative performance between these methods can depend on the loca-
 693 tion of interest and the aseasonal variance, but we believe that the machine learning method-
 694 ology can be widely applicable and competitive with optimal interpolation approaches
 695 globally. In particular, the Argo mixed layer depth samples with increased variance in
 696 the equatorial Pacific Ocean, whose subannual variability includes a relatively strong aseas-
 697 onal component, seem to be more strongly connected with the surface dynamics. There-
 698 fore, including surface information together with in-situ MLD estimates may be useful
 699 for generating improved reanalyses of the upper ocean under these circumstances. The
 700 second goal of this work was to use sophisticated probabilistic learning approaches to
 701 better understand the probability distribution of the MLD. The probabilistic approaches
 702 capture uncertainty to a greater extent than the optimal interpolation approach, but it
 703 is clear that, whether because of data or model limitations, more work is needed to ob-
 704 tain truly calibrated posterior probabilities. While initial results suggest that a Gaus-
 705 sian approximation of the conditional posterior distribution is appropriate, insufficient
 706 data might also explain the relative under-performance of the sampling-based probabilis-
 707 tic machine learning methods that we tested.

708 This work is an initial step into machine learning modeling of the MLD and there
 709 are several avenues for continued methodological and oceanographic research. First, the
 710 results in this study are regional test cases chosen to reveal how the variability of the
 711 MLD impacts the ability of the machine learning methods to learn a functional relation-
 712 ship between the surface variables and the MLD. Future work will expand this regional
 713 approach to a global scale. Second, while the probabilistic calibration results suggest that
 714 machine learning methods can better estimate the posterior distribution compared to
 715 the optimal interpolation approach, the overall calibration is underwhelming. Further
 716 research is needed to derive better architectures to better estimate this conditional pos-
 717 terior probability distribution. This research could include weight uncertainty, more so-
 718 phisticated sampling strategies, covariance regularization, or other neural network ar-
 719 chitectures. Finally, the research presented in this paper ignored temporal dynamics. We
 720 believe that incorporation of the temporal dynamics could help regularize the estima-
 721 tion procedure by coupling observations across time while simultaneously providing use-

ful scientific information about the temporal dynamics of the MLD in relation to the surface variables. In addition to the continued methodological research that follows from this paper, we believe that this methodology can be used to answer scientific oceanographic research questions that require fine resolution gridded MLD estimates.

Acknowledgments

This material is based upon work supported by the National Center for Atmospheric Research, which is a major facility sponsored by the NSF under Cooperative Agreement No. 1852977. Computing and data storage resources, including the Cheyenne supercomputer (doi:10.5065/D6RX99HX), were provided by the Computational and Information Systems Laboratory (CISL) at NCAR. DBW acknowledges useful and motivating discussions with NASA Salinity Science Team colleagues Justin Small, Ivana Cerovecki and Matt Mazloff on contract 80NSSC20K0890.

Code and examples for this project can be found at <https://github.com/NCAR/ml-ocean-bl> and <https://doi.org/10.5281/zenodo.4441098>. Argo-based mixed layer depth data (D. Whitt et al., 2020) can be accessed at <https://doi.org/10.5281/zenodo.4291175>. Preprocessed surface and mixed layer data and model outputs (Foster et al., 2020) can be accessed at <https://www.doi.org/10.5281/zenodo.4421752>. HYCOM data are obtained from <https://www.hycom.org/dataserver/gofs-3pt0/analysis>, experiment GLBu0.08 91.1 for 2014 and experiment GLBu0.08 19.1 for 2012.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., . . . Zheng, X. (2016). TensorFlow: A system for large-scale machine learning. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2016* (pp. 265–283). Retrieved from <https://tensorflow.org>.
- Ashukha, A., Lyzhov, A., Molchanov, D., & Vetrov, D. (2020, February). Pitfalls of In-Domain Uncertainty Estimation and Ensembling in Deep Learning. *arXiv e-prints*, arXiv:2002.06470.
- Ba, L. J., & Frey, B. (2013). Adaptive dropout for training deep neural networks. In *Advances in Neural Information Processing Systems*.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wierstra, D. (2015, 07–09 Jul). Weight uncertainty in neural network. In *Proceedings of the 32nd International Conference on Machine Learning* (Vol. 37, pp. 1613–1622). Lille, France: PMLR. Retrieved from <http://proceedings.mlr.press/v37/blundell115.html>
- Bolton, T., & Zanna, L. (2019, Jan). Applications of Deep Learning to Ocean Data Inference and Subgrid Parameterization. *Journal of Advances in Modeling Earth Systems*, 11(1), 376–399. Retrieved from <http://doi.wiley.com/10.1029/2018MS001472> doi: 10.1029/2018MS001472
- Brenowitz, N. D., & Bretherton, C. S. (2018, Jun). Prognostic Validation of a Neural Network Unified Physics Parameterization. *Geophysical Research Letters*, 45(12), 6289–6298. Retrieved from <http://doi.wiley.com/10.1029/2018GL078510> doi: 10.1029/2018GL078510
- Brooks, S. (2011). *Handbook of Markov chain Monte Carlo*. CRC Press/Taylor & Francis. Retrieved from <https://www.crcpress.com/Handbook-of-Markov-Chain-Monte-Carlo/Brooks-Gelman-Jones-Meng/p/book/9781420079418>
- Cabanes, C., Grouazel, A., von Schuckmann, K., Hamon, M., Turpin, V., Coatanoan, C., . . . Le Traon, P.-Y. (2013, Jan). The CORA dataset: validation and diagnostics of in-situ ocean temperature and salinity measurements. *Ocean Science*, 9(1), 1–18. Retrieved from <http://marine.copernicus.eu/services-portfolio/access-to-products/> doi: 10.5194/os-9-1-2013
- Caldeira, J., & Nord, B. (2020, April). Deeply Uncertain: Comparing Methods

- 773 of Uncertainty Quantification in Deep Learning Algorithms. *arXiv e-prints*,
 774 arXiv:2004.10710.
- 775 Caruana, R., Lawrence, S., & Giles, L. (2001, Jan). Overfitting in neural nets: Back-
 776 propagation, conjugate gradient, and early stopping. In *Advances in Neural*
 777 *Information Processing Systems*. Retrieved from [http://papers.nips.cc/
 778 paper/1895-overfitting-in-neural-nets-backpropagation-conjugate
 779 -gradient-and-early-stopping.pdf](http://papers.nips.cc/paper/1895-overfitting-in-neural-nets-backpropagation-conjugate-gradient-and-early-stopping.pdf)
- 780 Chaudhuri, A. H., Ponte, R. M., & Forget, G. (2016, Apr). Impact of uncertain-
 781 ties in atmospheric boundary conditions on ocean model solutions. *Ocean Mod-
 782 elling*, *100*, 96–108. doi: 10.1016/j.ocemod.2016.02.003
- 783 Cintra, R., De Campos Velho, H., Anochi, J., & Cocke, S. (2016, mar). Data as-
 784 similation by artificial neural networks for the global FSU atmospheric model:
 785 Surface pressure. In *2015 Latin-America Congress on Computational Intelli-
 786 gence, LA-CCI 2015*. Institute of Electrical and Electronics Engineers Inc. doi:
 787 10.1109/LA-CCI.2015.7435937
- 788 Cintra, R. S., & Velho, H. F. d. C. (2018, Jul). Data Assimilation by Artificial Neu-
 789 ral Networks for an Atmospheric General Circulation Model. In *Advanced ap-
 790 plications for artificial neural networks*. Retrieved from [http://arxiv.org/
 791 abs/1407.4360](http://arxiv.org/abs/1407.4360) doi: 10.5772/intechopen.70791
- 792 Cressie, N. A. C. (1993). *Statistics for Spatial Data* (Revised Ed ed.) (No. 1).
 793 Hoboken, NJ, USA: John Wiley & Sons, Inc. Retrieved from [http://
 794 doi.wiley.com/10.1002/9781119115151](http://doi.wiley.com/10.1002/9781119115151) doi: 10.1002/9781119115151
- 795 Cummings, J. A. (2006, Jan). Operational multivariate ocean data assimilation.
 796 *Quarterly Journal of the Royal Meteorological Society*, *131*(613), 3583–3604.
 797 Retrieved from <https://doi.org/10.1256/qj.05.105>
- 798 Cummings, J. A., & Smedstad, O. M. (2013). Variational data assimilation for the
 799 global ocean. In *Data assimilation for atmospheric, oceanic and hydrologic ap-
 800 plications (vol. ii)* (Vol. 2, pp. 303–343). Springer Berlin Heidelberg. doi: 10
 801 .1007/978-3-642-35088-7_13
- 802 Cybenko, G. (1989, Dec). Approximation by superpositions of a sigmoidal function.
 803 *Mathematics of Control, Signals, and Systems*, *2*(4), 303–314. Retrieved from
 804 <https://link.springer.com/article/10.1007/BF02551274> doi: 10.1007/
 805 BF02551274
- 806 Dillon, J. V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., ...
 807 Saurous, R. A. (2017, November). TensorFlow Distributions. *arXiv e-prints*,
 808 arXiv:1711.10604.
- 809 Dormann, C. F. (2020, Apr). Calibration of probability predictions from machine-
 810 learning and statistical models. *Global Ecology and Biogeography*, *29*(4), 760–
 811 765. Retrieved from [https://onlinelibrary.wiley.com/doi/abs/10.1111/
 812 geb.13070](https://onlinelibrary.wiley.com/doi/abs/10.1111/geb.13070) doi: 10.1111/geb.13070
- 813 Foster, D., Gagne II, D. J., & Whitt, D. (2020, Dec). *Probabilistic Machine Learn-
 814 ing Estimation of Ocean Mixed Layer Depth from Dense Satellite and Sparse
 815 In-Situ Observations: Preprocessed Satellite and In-situ observation datasets*.
 816 Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.4421752> doi:
 817 10.5281/zenodo.4421752
- 818 Fox, D. N., Teague, W. J., Barron, C. N., Carnes, M. R., & Lee, C. M. (2002, Feb).
 819 The modular ocean data assimilation system (MODAS). *Journal of Atmo-
 820 spheric and Oceanic Technology*, *19*(2), 240–252. Retrieved from [http://
 821 journals.ametsoc.org/jtech/article-pdf/19/2/240/3312918/1520-0426](http://journals.ametsoc.org/jtech/article-pdf/19/2/240/3312918/1520-0426)
 822 doi: 10.1175/1520-0426(2002)019<0240:TMODAS>2.0.CO;2
- 823 Frankignoul, C., & Hasselmann, K. (1977, Aug). Stochastic climate models, Part II
 824 Application to sea-surface temperature anomalies and thermocline variability.
 825 *Tellus*, *29*(4), 289–305. Retrieved from [http://tellusa.net/index.php/
 826 tellusa/article/view/11362](http://tellusa.net/index.php/tellusa/article/view/11362) doi: 10.1111/j.2153-3490.1977.tb00740.x
- 827 Gagne, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020,

- 828 mar). Machine Learning for Stochastic Parameterization: Generative Adversarial
 829 Networks in the Lorenz '96 Model. *Journal of Advances in Modeling Earth*
 830 *Systems*, 12(3). Retrieved from [https://onlinelibrary.wiley.com/doi/](https://onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001896)
 831 [abs/10.1029/2019MS001896](https://onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001896) doi: 10.1029/2019MS001896
- 832 Gal, Y. (2016). *Uncertainty in Deep Learning* (Unpublished doctoral dissertation).
 833 University of Cambridge.
- 834 Gal, Y., & Ghahramani, Z. (2016). Dropout as a Bayesian approximation: Representing
 835 model uncertainty in deep learning. In *33rd International Conference*
 836 *on Machine Learning, ICML 2016* (Vol. 3, pp. 1651–1660). Retrieved from
 837 <http://yarín.co>.
- 838 Gal, Y., Hron, J., & Kendall, A. (2017, Dec). Concrete dropout. In *Advances in*
 839 *Neural Information Processing Systems* (pp. 3582–3591).
- 840 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B.
 841 (2013). *Bayesian data analysis, third edition*. Taylor & Francis. Retrieved from
 842 <https://books.google.com/books?id=ZXL6AQAQBAJ>
- 843 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018, Jun).
 844 Could Machine Learning Break the Convection Parameterization Deadlock?
 845 *Geophysical Research Letters*, 45(11), 5742–5751. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/full/10.1029/2018GL078202)
 846 agupubs.onlinelibrary.wiley.com/doi/full/10.1029/2018GL078202 doi:
 847 10.1029/2018GL078202
- 848 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. Re-
 849 trieved from <http://www.deeplearningbook.org>
- 850 Guinehut, S., Dhomp, A.-L. L., Larnicol, G., & Le Traon, P.-Y. Y. (2012,
 851 Oct). High resolution 3-D temperature and salinity fields derived from
 852 in situ and satellite observations. *Ocean Science*, 8(5), 845–857. Re-
 853 trieved from <https://os.copernicus.org/articles/8/845/2012/> doi:
 854 10.5194/os-8-845-2012
- 855 Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of mod-
 856 ern neural networks. In *34th International Conference on Machine Learning,*
 857 *ICML 2017* (Vol. 3, pp. 2130–2143).
- 858 Hanawa, K., & Talley, L. D. (2001). Mode waters. *Ocean Circulation and Cli-*
 859 *mate: Observing and Modeling the Global Ocean*, 373–386 (736pp). Re-
 860 trieved from [ftp://bslctb.nerc-bas.ac.uk/jbsall/Papers_CMIP5team/](ftp://bslctb.nerc-bas.ac.uk/jbsall/Papers_CMIP5team/2001Hanawa.pdf)
 861 [2001Hanawa.pdf](ftp://bslctb.nerc-bas.ac.uk/jbsall/Papers_CMIP5team/2001Hanawa.pdf)
- 862 Hernández-Lobato, J. M., & Adams, R. P. (2015). Probabilistic backpropagation
 863 for scalable learning of Bayesian neural networks. In *32nd International Con-*
 864 *ference on Machine Learning, ICML 2015* (Vol. 3, pp. 1861–1869).
- 865 Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. R.
 866 (2012, July). Improving neural networks by preventing co-adaptation of feature
 867 detectors. *arXiv e-prints*, arXiv:1207.0580.
- 868 Hoffman, M. D., & Blei, D. M. (2015). Structured stochastic variational inference.
 869 In *Journal of Machine Learning Research* (Vol. 38, pp. 361–369). Retrieved
 870 from <http://jmlr.org/papers/v14/hoffman13a.html>
- 871 Holte, J., Talley, L. D., Gilson, J., & Roemmich, D. (2017, Jun). An Argo mixed
 872 layer climatology and database. *Geophysical Research Letters*, 44(11), 5618–
 873 5626. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/](https://agupubs.onlinelibrary.wiley.com/doi/full/10.1002/2017GL073426)
 874 [full/10.1002/2017GL073426](https://agupubs.onlinelibrary.wiley.com/doi/full/10.1002/2017GL073426)[https://agupubs.onlinelibrary.wiley](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL073426)
 875 [.com/doi/abs/10.1002/2017GL073426](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL073426)[https://agupubs.onlinelibrary.wiley](https://agupubs.onlinelibrary.wiley.com/doi/10.1002/2017GL073426)
 876 [.com/doi/10.1002/2017GL073426](https://agupubs.onlinelibrary.wiley.com/doi/10.1002/2017GL073426) doi: 10.1002/2017GL073426
- 877 Hornik, K., Stinchcombe, M., & White, H. (1989, Jan). Multilayer feedforward net-
 878 works are universal approximators. *Neural Networks*, 2(5), 359–366. doi: 10
 879 .1016/0893-6080(89)90020-8
- 880 Hsieh, W. W., & Tang, B. (1998, Sep). Applying Neural Network Models to Pre-
 881 diction and Data Analysis in Meteorology and Oceanography. *Bulletin of the*
 882 *American Meteorological Society*, 79(9), 1855–1870. Retrieved from <http://>

- 883 journals.ametsoc.org/bams/article-pdf/79/9/1855/3731484/1520-0477
 884 doi: 10.1175/1520-0477(1998)079<1855:ANNMTP>2.0.CO;2
- 885 Irrgang, C., Saynisch-Wagner, J., & Thomas, M. (2020, May). Machine Learning-
 886 Based Prediction of Spatiotemporal Uncertainties in Global Wind Velocity
 887 Reanalyses. *Journal of Advances in Modeling Earth Systems*, 12(5). Retrieved
 888 from <https://onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001876>
 889 doi: 10.1029/2019MS001876
- 890 Jiang, G. Q., Xu, J., & Wei, J. (2018, Apr). A Deep Learning Algorithm of Neu-
 891 ral Network for the Parameterization of Typhoon-Ocean Feedback in Ty-
 892 phoon Forecast Models. *Geophysical Research Letters*, 45(8), 3706–3716. doi:
 893 10.1002/2018GL077004
- 894 Kingma, D. P., & Ba, J. L. (2015, Dec). Adam: A method for stochastic optimiza-
 895 tion. In *3rd International Conference on Learning Representations, ICLR 2015*
 896 - *Conference Track Proceedings*. International Conference on Learning Repre-
 897 sentations, ICLR. Retrieved from <https://arxiv.org/abs/1412.6980v9>
- 898 Kingma, D. P., Salimans, T., & Welling, M. (2015). Variational dropout and the
 899 local reparameterization trick. In *Advances in Neural Information Processing*
 900 *Systems* (pp. 2575–2583).
- 901 Kingma, D. P., & Welling, M. (2014, Dec). Auto-encoding variational bayes. In *2nd*
 902 *International Conference on Learning Representations, ICLR 2014 - Confer-*
 903 *ence Track Proceedings*. International Conference on Learning Representations,
 904 ICLR. Retrieved from <https://arxiv.org/abs/1312.6114v10>
- 905 Kraus, E. B., & Turner, J. S. (1967, Jan). A one-dimensional model of the sea-
 906 sonal thermocline II. The general theory and its consequences. *Tellus*, 19(1),
 907 98–106. Retrieved from [https://www.tandfonline.com/doi/abs/10.3402/](https://www.tandfonline.com/doi/abs/10.3402/tellusa.v19i1.9753)
 908 [tellusa.v19i1.9753](https://www.tandfonline.com/doi/abs/10.3402/tellusa.v19i1.9753) doi: 10.3402/tellusa.v19i1.9753
- 909 Kuleshov, V., Fenner, N., & Ermon, S. (2018, Jul). Accurate uncertainties for deep
 910 learning using calibrated regression. In *35th International Conference on Ma-*
 911 *chine Learning, ICML 2018* (Vol. 6, pp. 4369–4377). PMLR. Retrieved from
 912 <http://proceedings.mlr.press/v80/kuleshov18a.html>
- 913 Labach, A., Salehinejad, H., & Valaee, S. (2019, April). Survey of Dropout Methods
 914 for Deep Neural Networks. *arXiv e-prints*, arXiv:1904.13310.
- 915 Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017, Dec). Simple and scalable
 916 predictive uncertainty estimation using deep ensembles. In *Advances in Neural*
 917 *Information Processing Systems* (pp. 6403–6414).
- 918 Lary, D. J., Alavi, A. H., Gandomi, A. H., & Walker, A. L. (2016, Jan). Machine
 919 learning in geosciences and remote sensing. *Geoscience Frontiers*, 7(1), 3–10.
 920 doi: 10.1016/j.gsf.2015.07.003
- 921 Leshno, M., Lin, V. Y., Pinkus, A., & Schocken, S. (1993, Jan). Multilayer
 922 feedforward networks with a nonpolynomial activation function can ap-
 923 proximate any function. *Neural Networks*, 6(6), 861–867. doi: 10.1016/
 924 [S0893-6080\(05\)80131-5](https://doi.org/10.1016/S0893-6080(05)80131-5)
- 925 Maeda, S.-i. (2014, December). A Bayesian encourages dropout. *arXiv e-prints*,
 926 arXiv:1412.7003.
- 927 McGovern, A., Elmore, K. L., Gagne, D. J., Haupt, S. E., Karstens, C. D.,
 928 Lagerquist, R., . . . Williams, J. K. (2017, Oct). Using artificial intelli-
 929 gence to improve real-time decision-making for high-impact weather. *Bul-*
 930 *letin of the American Meteorological Society*, 98(10), 2073–2090. doi:
 931 10.1175/BAMS-D-16-0123.1
- 932 Melnichenko, O., Hacker, P., Maximenko, N., Lagerloef, G., & Potemra, J. (2016,
 933 Jan). Optimum interpolation analysis of Aquarius sea surface salinity.
 934 *Journal of Geophysical Research: Oceans*, 121(1), 602–615. Retrieved from
 935 <https://onlinelibrary.wiley.com/doi/abs/10.1002/2015JC011343> doi:
 936 10.1002/2015JC011343
- 937 Monteleoni, C., Schmidt, G. A., & McQuade, S. (2013, Sep). Climate informatics:

- 938 Accelerating discovering in climate science with machine learning. *Computing*
 939 *in Science and Engineering*, 15(5), 32–40. doi: 10.1109/MCSE.2013.50
- 940 Neal, R. (1996). Bayesian Learning for Neural Networks. *Lecture Notes in Statistics*,
 941 1(118).
- 942 Nixon, J., Dusenberry, M., Jerfel, G., Nguyen, T., Liu, J., Zhang, L., & Tran, D.
 943 (2019, April). Measuring Calibration in Deep Learning. *arXiv e-prints*,
 944 arXiv:1904.01685.
- 945 O’Gorman, P. A., & Dwyer, J. G. (2018). Using Machine Learning to Parameterize
 946 Moist Convection: Potential for Modeling of Climate, Climate Change, and
 947 Extreme Events. *Journal of Advances in Modeling Earth Systems*, 10(10),
 948 2548–2563. doi: 10.1029/2018MS001351
- 949 Ouali, D., Chebana, F., & Ouarda, T. B. (2017, Jun). Fully nonlinear statistical
 950 and machine-learning approaches for hydrological frequency estimation at
 951 ungauged sites. *Journal of Advances in Modeling Earth Systems*, 9(2), 1292–
 952 1306. Retrieved from [https://onlinelibrary.wiley.com/doi/abs/10.1002/](https://onlinelibrary.wiley.com/doi/abs/10.1002/2016MS000830)
 953 [2016MS000830](https://onlinelibrary.wiley.com/doi/abs/10.1002/2016MS000830) doi: 10.1002/2016MS000830
- 954 Paisley, J., Blei, D. M., & Jordan, M. I. (2012). Variational Bayesian inference with
 955 stochastic search. In *Proceedings of the 29th International Conference on Ma-*
 956 *chine Learning, ICML 2012* (Vol. 2, pp. 1367–1374).
- 957 Pathak, J., Hunt, B., Girvan, M., Lu, Z., & Ott, E. (2018, Jan). Model-Free
 958 Prediction of Large Spatiotemporally Chaotic Systems from Data: A Reser-
 959 voir Computing Approach. *Physical Review Letters*, 120(2), 024102.
 960 Retrieved from [https://journals.aps.org/prl/abstract/10.1103/](https://journals.aps.org/prl/abstract/10.1103/PhysRevLett.120.024102)
 961 [PhysRevLett.120.024102](https://journals.aps.org/prl/abstract/10.1103/PhysRevLett.120.024102) doi: 10.1103/PhysRevLett.120.024102
- 962 Rasmussen, C. E., & Williams, C. K. I. (2006). *Gaussian Processes for Machine*
 963 *Learning*. MIT Press. Retrieved from [https://mitpress.mit.edu/books/](https://mitpress.mit.edu/books/gaussian-processes-machine-learning)
 964 [gaussian-processes-machine-learning](https://mitpress.mit.edu/books/gaussian-processes-machine-learning)
- 965 Rasp, S., Pritchard, M. S., & Gentine, P. (2018, Sep). Deep learning to represent
 966 subgrid processes in climate models. *Proceedings of the National Academy of*
 967 *Sciences of the United States of America*, 115(39), 9684–9689. doi: 10.1073/
 968 [pnas.1810286115](https://doi.org/10.1073/pnas.1810286115)
- 969 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais,
 970 N., & Prabhat. (2019, Feb). Deep learning and process understanding
 971 for data-driven Earth system science. *Nature*, 566(7743), 195–204. Re-
 972 trieved from <https://www.nature.com/articles/s41586-019-0912-1> doi:
 973 [10.1038/s41586-019-0912-1](https://doi.org/10.1038/s41586-019-0912-1)
- 974 Remote Sensing Systems. (2017). *GHRSSST Level 4 MW_OI Global Foundation*
 975 *Sea Surface Temperature analysis version 5.0 from REMSS*. NASA Physical
 976 Oceanography DAAC. Retrieved from [https://podaac.jpl.nasa.gov/](https://podaac.jpl.nasa.gov/dataset/MW_OI-REMSS-L4-GLOB-v5.0)
 977 [dataset/MW_OI-REMSS-L4-GLOB-v5.0](https://podaac.jpl.nasa.gov/dataset/MW_OI-REMSS-L4-GLOB-v5.0) doi: [https://doi.org/10.5067/](https://doi.org/10.5067/GHMWO-4FR05)
 978 [GHMWO-4FR05](https://doi.org/10.5067/GHMWO-4FR05)
- 979 Rezende, D. J., Mohamed, S., & Wierstra, D. (2014, Jan). Stochastic backpropa-
 980 gation and approximate inference in deep generative models. In *31st Interna-*
 981 *tional Conference on Machine Learning, ICML 2014* (Vol. 4, pp. 3057–3070).
 982 Retrieved from <http://proceedings.mlr.press/v32/rezende14.html>
- 983 Roemmich, D., & Gilson, J. (2009, Aug). The 2004–2008 mean and annual
 984 cycle of temperature, salinity, and steric height in the global ocean from
 985 the Argo Program. *Progress in Oceanography*, 82(2), 81–100. doi:
 986 [10.1016/j.pocean.2009.03.004](https://doi.org/10.1016/j.pocean.2009.03.004)
- 987 Rosenblatt, F. (1958, Nov). The perceptron: A probabilistic model for information
 988 storage and organization in the brain. *Psychological Review*, 65(6), 386–408.
 989 doi: 10.1037/h0042519
- 990 Ruder, S. (2016, September). An overview of gradient descent optimization algo-
 991 rithms. *arXiv e-prints*, arXiv:1609.04747.
- 992 Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representa-

993 tions by back-propagating errors. *Nature*, 323(6088), 533–536. Retrieved from
 994 <https://www.nature.com/articles/323533a0> doi: 10.1038/323533a0

995 Schmidtko, S., Johnson, G. C., & Lyman, J. M. (2013, Apr). MIMOC: A global
 996 monthly isopycnal upper-ocean climatology with mixed layers. *Journal of*
 997 *Geophysical Research: Oceans*, 118(4), 1658–1672. Retrieved from [https://](https://agupubs.onlinelibrary.wiley.com/doi/full/10.1002/jgrc.20122)
 998 agupubs.onlinelibrary.wiley.com/doi/full/10.1002/jgrc.20122 doi:
 999 10.1002/jgrc.20122

1000 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R.
 1001 (2014). Dropout: A simple way to prevent neural networks from overfitting.
 1002 *Journal of Machine Learning Research*, 15, 1929–1958.

1003 Stommel, H. (1979, Jul). Determination of water mass properties of water
 1004 pumped down from the Ekman layer to the geostrophic flow below. *Pro-*
 1005 *ceedings of the National Academy of Sciences*, 76(7), 3051–3055. Retrieved
 1006 from <https://www.pnas.org/content/76/7/3051>[https://www.pnas.org/](https://www.pnas.org/content/76/7/3051.abstract)
 1007 [content/76/7/3051.abstract](https://www.pnas.org/content/76/7/3051.abstract) doi: 10.1073/pnas.76.7.3051

1008 Ukkonen, P., & Mäkelä, A. (2019, Jun). Evaluation of Machine Learning Classifiers
 1009 for Predicting Deep Convection. *Journal of Advances in Modeling Earth Sys-*
 1010 *tems*, 11(6), 1784–1802. Retrieved from [https://onlinelibrary.wiley.com/](https://onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001561)
 1011 [doi/abs/10.1029/2018MS001561](https://onlinelibrary.wiley.com/doi/abs/10.1029/2018MS001561) doi: 10.1029/2018MS001561

1012 Valler, V., Franke, J., & Brönnimann, S. (2019). Impact of different estimations
 1013 of the background-error covariance matrix on climate reconstructions based
 1014 on data assimilation. *Climate of the Past*, 15(4), 1427–1441. Retrieved from
 1015 <https://doi.org/10.5194/cp-15-1427-2019> doi: 10.5194/cp-15-1427-2019

1016 Wahle, K., Staneva, J., & Guenther, H. (2015, Dec). Data assimilation of ocean
 1017 wind waves using Neural Networks: A case study for the German Bight. *Ocean*
 1018 *Modelling*, 96, 117–125. doi: 10.1016/j.ocemod.2015.07.007

1019 Welling, M., & Teh, Y. W. (2011). Bayesian learning via stochastic gradient
 1020 langevin dynamics. In *Proceedings of the 28th International Conference on*
 1021 *Machine Learning, ICML 2011* (pp. 681–688).

1022 Weyn, J. A., Durran, D. R., & Caruana, R. (2019, Aug). Can Machines
 1023 Learn to Predict Weather? Using Deep Learning to Predict Gridded 500-
 1024 hPa Geopotential Height From Historical Weather Data. *Journal of Ad-*
 1025 *vances in Modeling Earth Systems*, 11(8), 2680–2693. Retrieved from
 1026 <https://onlinelibrary.wiley.com/doi/abs/10.1029/2019MS001705> doi:
 1027 10.1029/2019MS001705

1028 Whitt, D., Nicholson, S., & Carranza, M. (2020, November). *Argo-based ocean*
 1029 *surface mixed layer depths using the buoyancy gradient definition of Whitt*
 1030 *Nicholson and Carranza (2019)*. Zenodo. Retrieved from [https://doi.org/](https://doi.org/10.5281/zenodo.4291175)
 1031 [10.5281/zenodo.4291175](https://doi.org/10.5281/zenodo.4291175) doi: 10.5281/zenodo.4291175

1032 Whitt, D. B., Nicholson, S. A., & Carranza, M. M. (2019, Dec). Global Impacts
 1033 of Subseasonal (<60 Day) Wind Variability on Ocean Surface Stress, Buoy-
 1034 ancy Flux, and Mixed Layer Depth. *Journal of Geophysical Research: Oceans*,
 1035 124(12), 8798–8831. Retrieved from [https://onlinelibrary.wiley.com/](https://onlinelibrary.wiley.com/doi/abs/10.1029/2019JC015166)
 1036 [doi/abs/10.1029/2019JC015166](https://onlinelibrary.wiley.com/doi/abs/10.1029/2019JC015166) doi: 10.1029/2019JC015166

1037 Zlotnicki, V., Qu, Z., & Willis, J. (2019). *MEaSUREs Gridded Sea Surface Height*
 1038 *Anomalies Version 1812*. NASA Physical Oceanography DAAC. Retrieved
 1039 from https://podaac.jpl.nasa.gov/dataset/SEA_SURFACE_HEIGHT_ALT
 1040 [_GRIDS_L4_2SATS_5DAY_6THDEG_V_JPL1812](https://podaac.jpl.nasa.gov/dataset/SEA_SURFACE_HEIGHT_ALT) doi: 10.5067/SLREF-CDRV2